
Kernel Methods and Introduction to Gaussian Processes

*Prof. Nicholas Zabaras
University of Notre Dame
Notre Dame, IN, USA*

*Email: nzabaras@gmail.com
URL: <https://www.zabaras.com/>*

November 13, 2017



Contents

- Keeping or discarding the training data, Kernel Function, Dual Representation, Constructing Kernels from Basis Functions, Constructing Kernels Directly
- Combining Kernels, Gaussian Kernel, Symbolic Input, Probabilistic Generative Models, Fisher Kernel, Sigmoidal Kernel
- Radial Basis Functions, Reducing the Size of the Basis, Parzen Density Estimator, Nadaraya-Watson model
- Gaussian Processes, Gaussian Process Regression, Gaussian process vs Basis Function Regression, Learning parameters
- Automatic relevance detection
- Gaussian Process Classification, Laplace Approximation, Connections to Bayesian Neural Networks

Following closely: Bishop CM, Pattern Recognition and Machine Learning, Springer, 2006 (Chapter 6)



Training Data: Keep or Discard?

- We have considered linear parametric methods e.g. for regression.
- The form of the mapping $y(\mathbf{x}, \mathbf{w})$ from input \mathbf{x} to output y is governed by \mathbf{w} .
- During the learning phase, the training data is used either
 - to obtain a point estimate of \mathbf{w} or
 - posterior distribution $p(\mathbf{w}|\mathbf{x}, \mathbf{t})$.
- The training data is then discarded, and predictions for new inputs are based purely on the learned parameter vector \mathbf{w} .
- The same approach is used in nonlinear parametric models including neural networks.

Training Data: Keep or Discard?

- There are machine learning techniques, in which the training data points are kept and used for prediction.
 - Parzen probability density model: set of kernel functions centered on training data points
 - Nearest neighbors technique: assigning to each new test vector the same label as the closest example from the training set
- These are **memory-based methods**: storing the training set in order to make predictions for future data points.
- They require a metric to measure the similarity of any two vectors in input space.
- They are generally fast to ‘train’ but slow at making predictions for test data points.

Kernel Methods and Kernel Trick

- Many linear parametric models can be recast into an equivalent ‘dual representation’
 - **Prediction is based on linear combinations of a kernel function evaluated at the training data points**
- For models which are based on a fixed nonlinear feature space mapping $\phi(x)$, the kernel function is given by

$$k(x, x') = \phi(x)^T \phi(x')$$

- Using the identity mapping for the feature space we obtain the **linear kernel**:

$$\phi(x) = x, \quad k(x, x') = x^T x'$$

- **Kernel trick:** if we have an algorithm formulated such that x enters only in the form of scalar products, then we can replace that scalar product with some other choice of kernel.

- Aizerman, M. A., E. M. Braverman, and L. I. Rozonoer (1964). [The probability problem of pattern recognition learning and the method of potential functions. Automation and Remote Control](#) **25**, 1175–1190.
- Boser, B. E., I. M. Guyon, and V. N. Vapnik (1992). [A training algorithm for optimal margin classifiers](#). In D. Haussler (Ed.), [Proceedings Fifth Annual Workshop on Computational Learning Theory](#) (COLT), pp. 144–152. ACM.Scholkopf et al.,1998).



Kernel Function

- For models based on feature space mapping $\phi(x)$:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

- Its a symmetric function:

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$$

- Linear kernel:

$$\phi(\mathbf{x}) = \mathbf{x}, \quad k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- Stationary kernel: invariant to translations in \mathbf{x}

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$$

- Homogeneous kernel (radial basis functions):

$$k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$$



Kernel Substitution Applications

- Algorithms expressed in terms of scalar products can be reformulated using kernel substitution (kernel trick). Applications include:

- Non-linear PCA,
- Nearest-neighbor classifiers,
- Kernel Fisher discriminant, etc.

- Schölkopf, B., A. Smola, and K.-R. Müller (1998). [Nonlinear component analysis as a kernel eigenvalue problem](#). *Neural Computation* **10**(5), 1299–1319.
- Mika, S., G. Rätsch, J. Weston, and B. Schölkopf (1999). [Fisher discriminant analysis with kernels](#). In Y. H. Hu, J. Larsen, E. Wilson, and S. Douglas (Eds.), [Neural Networks for Signal Processing IX](#), pp. 41–48. IEEE.
- Roth, V. and V. Steinlage (2000). [Nonlinear discriminant analysis using kernel functions](#). In S. A. Solla, T. K. Leen, and K. R. Müller (Eds.), [Advances in Neural Information Processing Systems](#), Volume 12. MIT Press.
- Baudat, G. and F. Anouar (2000). [Generalized discriminant analysis using a kernel approach](#). *Neural Computation* **12**(10), 2385–2404.

Textbooks on Kernel Methods include:

- Schölkopf, B. and A. J. Smola (2002). [Learning with Kernels](#). MIT Press.
- Herbrich, R. (2002). [Learning Kernel Classifiers](#). MIT Press.
- Shawe-Taylor, J. and N. Cristianini (2004). [Kernel Methods for Pattern Analysis](#). Cambridge University Press ([slides](#))



Dual Representation

- Many linear models for regression and classification can be reformulated in terms of a dual representation in which the kernel function arises naturally.
- Consider a linear regression model with a regularized cost functional ($\lambda > 0$).

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Set the gradient to zero:

$$\begin{aligned}\mathbf{w} &= -\frac{1}{\lambda} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a} \\ a_n &\equiv -\frac{1}{\lambda} (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n), \quad \mathbf{a} = \{a_1, \dots, a_N\}^T\end{aligned}$$

- Note that we defined \mathbf{a} using only the 1st term in $J(\mathbf{w})$.

Dual Representation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}(x_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}, \quad \mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}(x_n) - t_n) \boldsymbol{\phi}(x_n) = \Phi^T \mathbf{a}$$

□ Substitute \mathbf{w} in the cost functional:

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2} \sum_{n=1}^N (\mathbf{a}^T \Phi \boldsymbol{\phi}(x_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \Phi \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \end{aligned}$$

□ Here we used from an earlier lecture

$$\Phi = \begin{pmatrix} \boldsymbol{\phi}^T(x_1) \\ \boldsymbol{\phi}^T(x_2) \\ \vdots \\ \boldsymbol{\phi}^T(x_N) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\phi}_0(x_1) & \boldsymbol{\phi}_1(x_1) & .. & \boldsymbol{\phi}_{M-1}(x_1) \\ \boldsymbol{\phi}_0(x_2) & \boldsymbol{\phi}_1(x_2) & .. & \boldsymbol{\phi}_{M-1}(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ \boldsymbol{\phi}_0(x_N) & \boldsymbol{\phi}_1(x_N) & .. & \boldsymbol{\phi}_{M-1}(x_N) \end{pmatrix},$$

$$\begin{aligned} \Phi^T \Phi &= \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T, \\ \Phi^T \mathbf{t} &= \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n) \end{aligned}$$

Dual Representation: Gramm Matrix

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

- We can now define the **Gramm matrix (symmetric)** $\mathbf{K}_{N \times N}$ using the kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$ as follows:

$$K_{nm} = K_{mn} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) \quad \text{or} \quad \mathbf{K} = \Phi \Phi^T$$

$$\mathbf{K} = \Phi \Phi^T = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{pmatrix} (\phi(\mathbf{x}_1) \quad \phi(\mathbf{x}_2) \quad \dots \quad \phi(\mathbf{x}_N))$$

- The cost function takes the **dual representation** form:

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

Dual Representation

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

- Setting the gradient wrt \mathbf{a} equal to zero:

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_{N \times N})^{-1} \mathbf{t}, \quad K_{nm} = K_{mn} = k(\mathbf{x}_n, \mathbf{x}_m) \quad \text{or} \quad \mathbf{K} = \Phi \Phi^T \quad (N \times N \text{ matrix})$$

- Using $\mathbf{w} = \Phi^T \mathbf{a}$, the prediction for a new input \mathbf{x} is

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \phi(\mathbf{x})^T \Phi^T \mathbf{a} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_{N \times N})^{-1} \mathbf{t}$$

where

$$\begin{aligned} \mathbf{k}(\mathbf{x})^T &= \phi(\mathbf{x})^T \Phi^T = (\phi(\mathbf{x})^T \phi(\mathbf{x}_1) \quad \phi(\mathbf{x})^T \phi(\mathbf{x}_2) \quad \dots \quad \phi(\mathbf{x})^T \phi(\mathbf{x}_N)) \\ &= (k(\mathbf{x}_1, \mathbf{x}) \quad k(\mathbf{x}_2, \mathbf{x}) \quad \dots \quad k(\mathbf{x}_N, \mathbf{x})) \end{aligned}$$

- We now invert an $N \times N$ matrix instead of $M \times M$ (to compute \mathbf{w}). Everything is in terms of kernels (many choices here) but often $N \gg M$.

Summary of the Dual Representation

- The dual formulation allows the solution to the problem to be expressed entirely in terms of the kernel function $k(\mathbf{x}, \mathbf{x}')$.
- The prediction at \mathbf{x} is given by a linear combination of the target values from the training set.

$$y(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_{N \times N})^{-1} \mathbf{t}$$
$$\mathbf{k}(\mathbf{x})^T = (k(\mathbf{x}_1, \mathbf{x}) \ k(\mathbf{x}_2, \mathbf{x}) \ \dots \ k(\mathbf{x}_N, \mathbf{x}))$$

- We already have seen this result in a different form in the linear regression set of notes.
- Because we now have to invert a larger matrix $N \times N$, the formulation does not seem practical.
- However, we now have the choice of working with kernel functions rather than the basis functions $\phi(\mathbf{x})$. The future space can even be infinite dimensional.

Dual of the Dual Representation

- It is not difficult to see that the dual of the dual formulation is given by our original least squares representation for \mathbf{w} .
- Note that $J(\mathbf{a})$ can be written in terms of $\mathbf{K}\mathbf{a}$, $\mathbf{K} = \Phi\Phi^T$

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2}\mathbf{a}^T\Phi\Phi^T\Phi\Phi^T\mathbf{a} - \mathbf{a}^T\Phi\Phi^T\mathbf{t} + \frac{1}{2}\mathbf{t}^T\mathbf{t} + \frac{\lambda}{2}\mathbf{a}^T\Phi\Phi^T\mathbf{a} = \\ &= \frac{1}{2}(\mathbf{K}\mathbf{a} - \mathbf{t})^T(\mathbf{K}\mathbf{a} - \mathbf{t}) + \frac{\lambda}{2}\mathbf{a}^T\mathbf{K}\mathbf{a} \end{aligned}$$

- Recall that:

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} \phi^T(x_1) \\ \phi^T(x_2) \\ \vdots \\ \phi^T(x_N) \end{pmatrix} (\phi(x_1) \quad \phi(x_2) \quad .. \quad \phi(x_N))$$

- Since $N \gg M$, \mathbf{K} is rank deficient – $N - M$ eigenvectors of \mathbf{K} have zero eigenvalues. We can decompose \mathbf{a} as: $\mathbf{a} = \mathbf{a}_{\parallel} + \mathbf{a}_{\perp}$, $\mathbf{a}_{\perp}^T\mathbf{a}_{\parallel} = 0$ and $\mathbf{K}\mathbf{a}_{\perp} = \mathbf{0}$. Thus \mathbf{a}_{\perp} is not determined by $J(\mathbf{a})$. One can set $\mathbf{a}_{\perp} = \mathbf{0}$, and thus

$$\mathbf{a} = \mathbf{a}_{\parallel} \in \text{span}\mathbf{K} \quad (\mathbf{K} = \Phi\Phi^T) \Rightarrow \mathbf{a} = \Phi(\Phi^T\mathbf{v}) \Rightarrow \mathbf{a} = \Phi\mathbf{u}$$

Dual of the Dual Representation

- The cost functional becomes:

$$J(\mathbf{u}) = \frac{1}{2} (\Phi \Phi^T \Phi \mathbf{u} - \mathbf{t})^T (\Phi \Phi^T \Phi \mathbf{u} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{u}^T \Phi^T \Phi \Phi^T \Phi \mathbf{u}$$

- Since $\Phi^T \Phi$ has full rank, we can introduce an equivalent parametrization as

$$\mathbf{w} = \Phi^T \Phi \mathbf{u}$$

- The cost functional now becomes:

$$J(\mathbf{u}) = \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w}$$

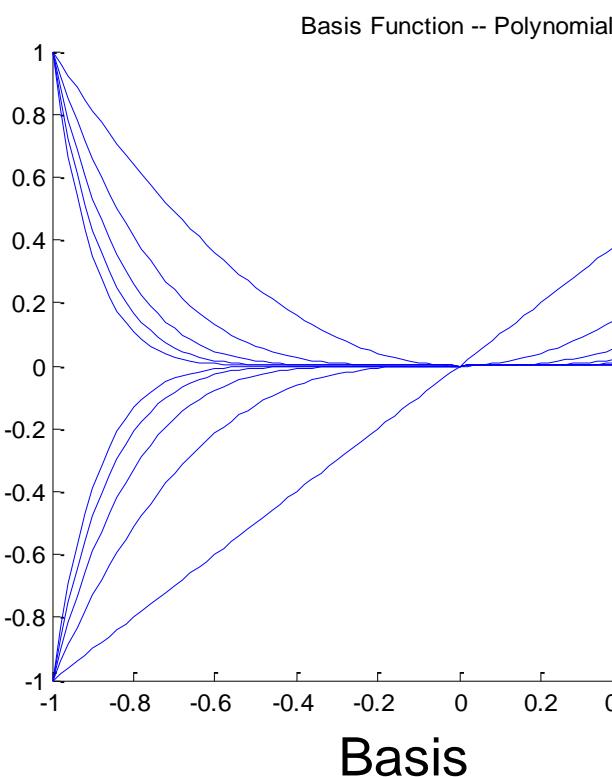
- This is our original least squares problem. So the proof is complete.

Constructing Kernels From Basis Functions

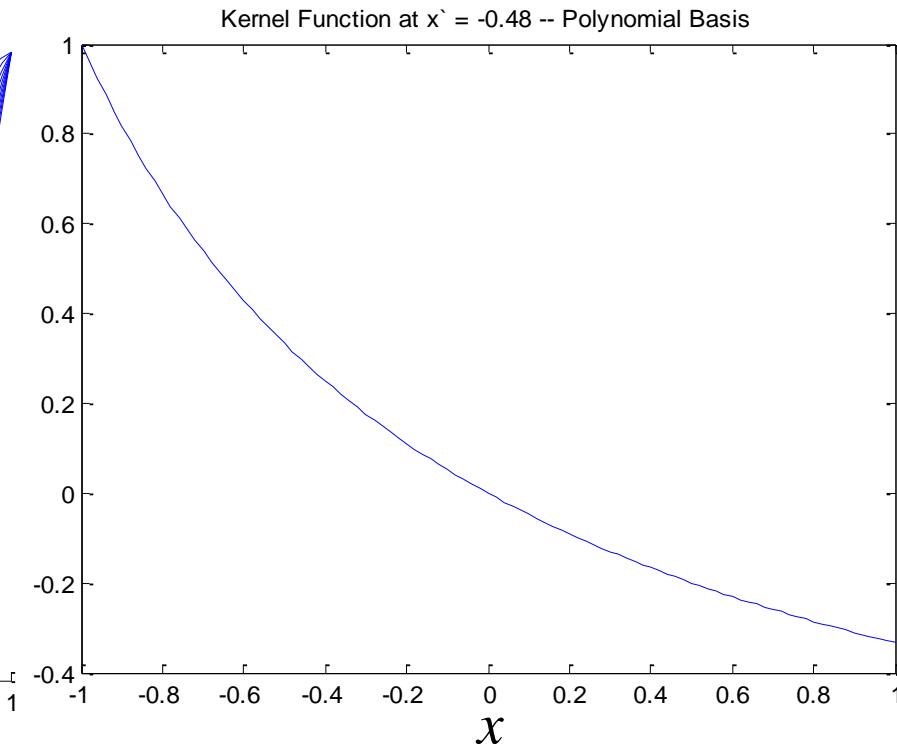
- Choose feature space mapping $\phi(x)$, then the kernel (in 1D) is given as:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

Polynomials



MatLab Code

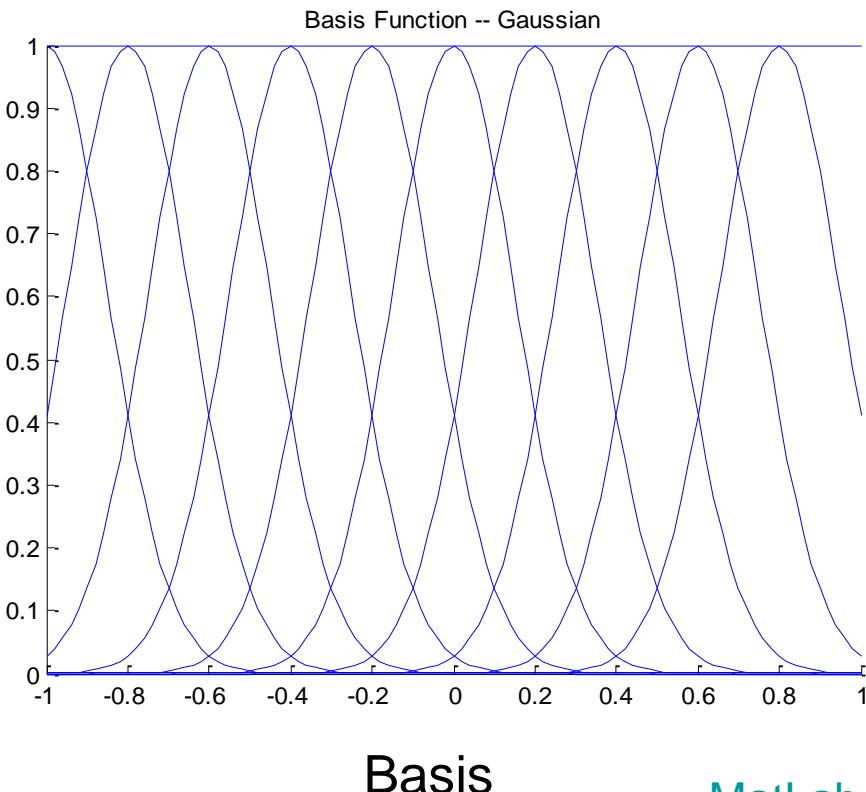


Constructing Kernels From Basis Functions

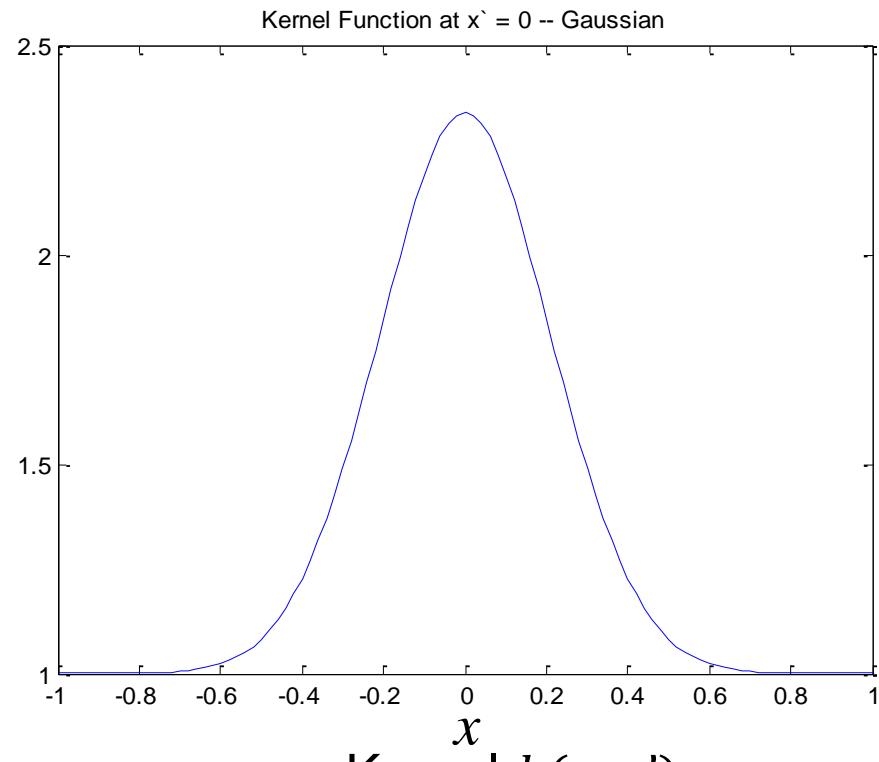
- Choose feature space mapping $\phi(x)$, then the kernel (in 1D) is given as:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

Gaussians



MatLab Code

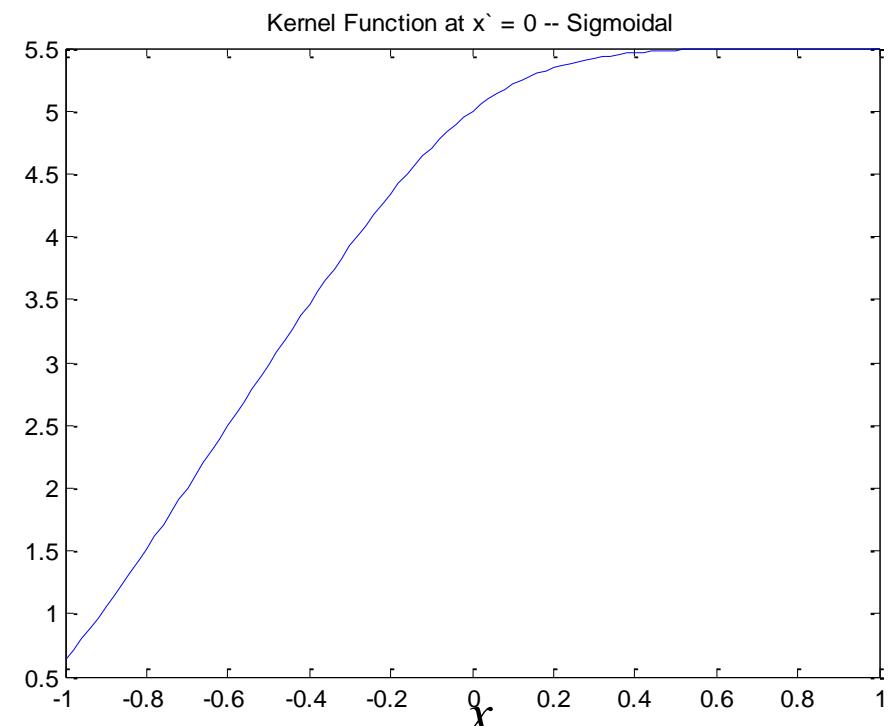
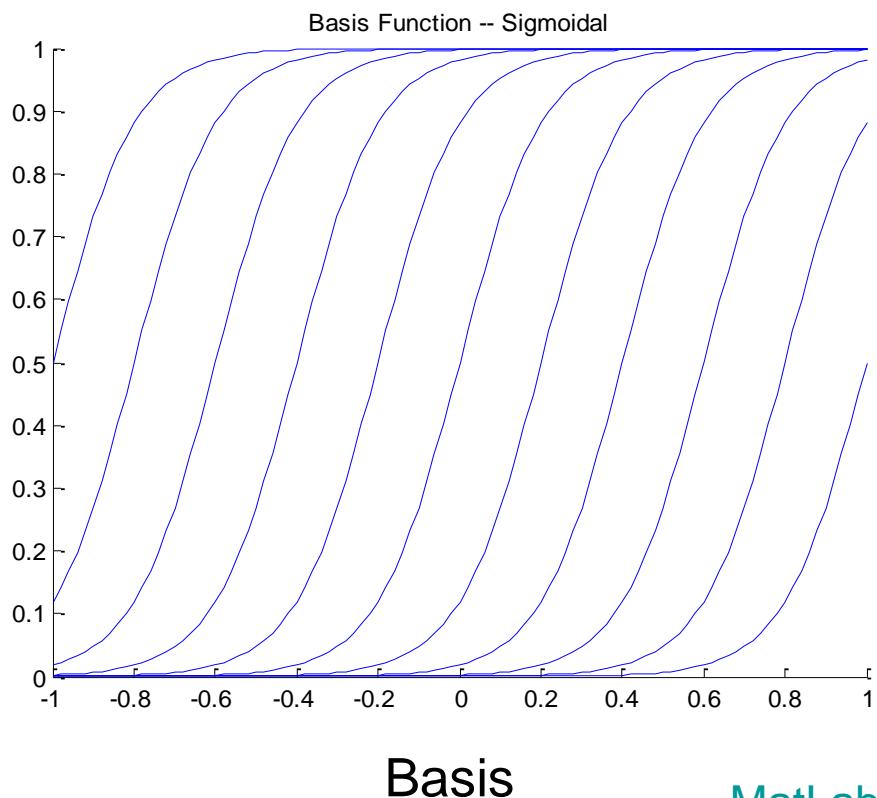


Constructing Kernels From Basis Functions

- Choose feature space mapping $\phi(x)$, then the kernel (in 1D) is given as:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

Logistic Sigmoid



MatLab Code



Constructing Kernels Directly

- When constructing the kernel functional directly verify its validity
 - it should correspond to a scalar product in some feature space

- Simple example

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

- In the 2-D case this corresponds to

$$k(\mathbf{x}, \mathbf{z}) = (x_1 z_1 + x_2 z_2)^2 = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

where

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$$

- To test validity without having to construct $\phi(\mathbf{x})$ explicitly, one can check if $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$ is semidefinite:

Function $k(\mathbf{x}_n, \mathbf{x}_m)$ is a valid kernel $\Leftrightarrow K \geq 0$ (Gram matrix) for all $\{\mathbf{x}_n\}$

- Shawe-Taylor, J. and N. Cristianini (2004). *Kernel Methods for Pattern Analysis*. Cambridge University Press ([slides](#))



Kernel Engineering

- Build kernels using simple kernels as building blocks. Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$, $k_2(\mathbf{x}, \mathbf{x}')$, the following kernels are valid:

$$\begin{array}{ll} k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}'), & k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \\ k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')), & k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \\ k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}'), & k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \\ k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')), & k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T A \mathbf{x}' \\ k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b), & k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \end{array}$$

with corresponding conditions:

$c > 0$, $f(\cdot)$ arbitrary, $q(\cdot)$ polynomial with ≥ 0 coefficients,

$\phi(\mathbf{x}) \in \mathbb{R}^M$, $k_3(\cdot, \cdot)$ valid kernel in \mathbb{R}^M , A sym. ≥ 0 ,

$\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, k_a , k_b valid kernels

Combining Kernels: Kernel Engineering

- If $k_1(x, x')$ is a valid kernel, so is $k(x, x') = ck_1(x, x')$, $c > 0$

$$k(x, x') = ck_1(x, x') = c\phi(x)^T \phi(x') = u(x)^T u(x'), \text{ where } u(x) = \sqrt{c}\phi(x)$$

- If $k_1(x, x')$ is a valid kernel so is $k(x, x') = f(x)k_1(x, x')f(x')$ for $f(\cdot)$ arbitrary.

$$k(x, x') = f(x)k_1(x, x')f(x') = f(x)\phi(x)^T \phi(x')f(x') = u(x)^T u(x'), \text{ where } u(x) = f(x)\phi(x)$$

- If $k_1(x, x'), k_2(x, x')$ are valid kernels so is $k_1(x, x') + k_2(x, x')$.
Indeed, note that the Gram matrix K with elements $k(x_n, x_m)$ for this kernel is ≥ 0 .

$$\forall a, a^T K a = a^T K_1 a + a^T K_2 a \geq 0, \text{ since } K_1, K_2 \geq 0$$



Kernel Examples

□ Polynomial kernels:

- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$: $\phi(\mathbf{x})$ contains only terms of degree 2
- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^2, c > 0$: $\phi(\mathbf{x})$ contains constant, linear and order 2 terms
- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^M$: contains all monomials of order M

Example: For \mathbf{x} and \mathbf{x}' images, $k(\mathbf{x}, \mathbf{x}')$ is a weighted sum of all possible products of M pixels of image \mathbf{x} with M pixels of image \mathbf{x}' .

- $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M, c > 0$: contains all terms up to order M



Gaussian Kernel

- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

- To see that this is a valid kernel note:

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'$$

and write it as:

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}')\right) \\ &= \underbrace{\exp\left(-\frac{\mathbf{x}^T \mathbf{x}}{2\sigma^2}\right)}_{f(\mathbf{x})} \underbrace{\exp\left(\frac{\mathbf{x}^T \mathbf{x}'}{\sigma^2}\right)}_{k_1(\mathbf{x}, \mathbf{x}')} \underbrace{\exp\left(-\frac{\mathbf{x}'^T \mathbf{x}'}{2\sigma^2}\right)}_{f(\mathbf{x}')} \end{aligned}$$

- “Taylor expanding” the middle exponential, we can see that the Gaussian kernel can be expressed as the inner product of an infinite dimensional feature vector.



Gaussian Kernel: non-Euclidean Distance

- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}')\right)$$

- With **the kernel trick**, we can replace the Euclidean distance with any non-linear kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2}(\kappa(\mathbf{x}, \mathbf{x}) + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}'))\right)$$



Example: Symbolic Inputs

- The inputs can be symbolic, graphs, sets, strings, text, etc.
- For example, fix a set D and consider all of its subsets ($2^{|D|}$). If A_1 and A_2 are two subsets of D , we can define the kernel as:

$$k(A_1, A_2) = 2^{|A_1 \cap A_2|}$$

where $|A_1 \cap A_2|$ is the number of elements in the intersection $A_1 \cap A_2$

- $\phi(D)$ is defined to map to a vector of $2^{|D|}$ 1's, one for each possible subset of D , including D itself as well as the empty set.
- For $A \subset D$, $\phi(A)$ will have 1's in all positions that correspond to subsets of A and 0's in all other positions. Therefore, $\phi(A_1)^T \phi(A_2)$ will count the number of subsets shared by A_1 and A_2 .
- This can also be obtained by counting the number of elements in the intersection of A_1 and A_2 , and then raising 2 to this number, which is exactly $k(A_1, A_2) = 2^{|A_1 \cap A_2|}$

Probabilistic Generative Models

- Use a generative model to define a kernel and then use this kernel in a discriminative approach.
- Kernel for generative model $p(x)$:

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}')$$

- This is a valid kernel function since it is an inner product in the 1D feature space defined by the mapping $p(\mathbf{x})$.
- It says that two inputs \mathbf{x} and \mathbf{x}' are similar if they both have high probabilities.

- Haussler, D. (1999). [Convolution kernels on discrete structures](#). Technical Report UCSC-CRL-99-10, University of California, Santa Cruz, Computer Science Department.
- Lasserre, J., C. M. Bishop, and T. Minka (2006). [Principled hybrids of generative and discriminative models](#). In *Proceedings 2006 IEEE Conference on Computer Vision and Pattern Recognition*, New York.



Probabilistic Generative Models

- Using $k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}')$ and the kernel engineering formulas

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}'), c > 0$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

we can generalize as:

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \sum_i p(\mathbf{x} | i)p(\mathbf{x}' | i)p(i) \\ k(\mathbf{x}, \mathbf{x}') &= \int p(\mathbf{x} | \mathbf{z})p(\mathbf{x}' | \mathbf{z})p(\mathbf{z})d\mathbf{z} \end{aligned}$$

i and \mathbf{z} are latent variables.

- \mathbf{x} and \mathbf{x}' will give large value of $k(\mathbf{x}, \mathbf{x}')$ if they have significant probability under a range of components.



Kernel Function in HMMs

- Now consider an observation is given by a sequence $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$.
- A popular generative model for sequences is the **Hidden Markov model (HMM)**, which expresses the distribution $p(\mathbf{X})$ as a marginalization over a corresponding sequence of hidden states $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_L\}$.
- We define **a kernel function measuring the similarity of two sequences \mathbf{X} and \mathbf{X}'** by extending the earlier mixture representation

$$k(\mathbf{X}, \mathbf{X}') = \sum_{\mathbf{Z}} p(\mathbf{X} | \mathbf{Z}) p(\mathbf{X}' | \mathbf{Z}) p(\mathbf{Z})$$

- Note that **both observed sequences \mathbf{X}, \mathbf{X}' are generated by the same hidden sequence \mathbf{Z}** .



Fisher Kernel

- Consider a parametric generative model $p(\mathbf{x}|\boldsymbol{\theta})$. We want a kernel that measures the similarity of two input vectors \mathbf{x} and \mathbf{x}' induced by this distribution.
- The Fisher score is introduced defining a vector in future space

$$g(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x} | \boldsymbol{\theta})$$

- The Fisher kernel is defined as

$$k(\mathbf{x}, \mathbf{x}') = g(\boldsymbol{\theta}, \mathbf{x})^T \mathbf{F}^{-1} g(\boldsymbol{\theta}, \mathbf{x}')$$

The Fisher information matrix \mathbf{F} (considers the geometry of the parameter/future space) is given as the expectation under $p(\mathbf{x} | \boldsymbol{\theta})$

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} \left[g(\boldsymbol{\theta}, \mathbf{x}) g(\boldsymbol{\theta}, \mathbf{x})^T \right]$$

- Using \mathbf{F} makes the kernel invariant under $\boldsymbol{\theta} \rightarrow \psi(\boldsymbol{\theta})$ where $\psi(\boldsymbol{\theta})$ is non-linear reparametrization that is invertible & differentiable.

- Jaakkola, T. S. and D. Haussler (1999). [Exploiting generative models in discriminative classifiers](#). In M. S. Kearns, S. A. Solla, and D. A. Cohn (Eds.), [Advances in Neural Information Processing Systems](#), Volume 11. MIT Press.
- Amari, S. I. (1998). [Natural gradient works efficiently in learning](#). [Neural Computation 10, 251–276](#).

Fisher Kernel

$$\mathbf{F} = \mathbb{E}_x \left[g(\theta, x) g(\theta, x)^T \right]$$

- We often approximate the information matrix with a sample average:

$$\mathbf{F} \approx \frac{1}{N} \sum_{n=1}^N g(\theta, \mathbf{x}_n) g(\theta, \mathbf{x}_n)^T$$

- This is the covariance matrix of the Fisher scores, and so the Fisher kernel corresponds to a whitening of these scores.

$$k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x})^T \mathbf{F}^{-1} g(\theta, \mathbf{x}')$$

- More simply, we can just omit the Fisher information matrix altogether and use the non-invariant kernel

$$k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x})^T g(\theta, \mathbf{x}')$$

▪ Hofmann, T. (2000). [Learning the similarity of documents: an information-geometric approach to document retrieval and classification](#). In S. A. Solla, T. K. Leen, and K. R. Müller (Eds.), [Advances in Neural Information Processing Systems](#), Volume 12, pp. 914–920. MIT Press.



Fisher Kernel: Example

- Let us compute the Fisher kernel for the Gaussian distribution with fixed covariance S :

$$p(\mathbf{x} | \boldsymbol{\mu}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, S)$$

- The Fisher score is given:

$$g(\boldsymbol{\mu}, \mathbf{x}) = \nabla_{\boldsymbol{\mu}} \ln p(\mathbf{x} | \boldsymbol{\mu}) = S^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

- The Fisher information matrix \mathbf{F} is then:

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} \left[g(\boldsymbol{\mu}, \mathbf{x}) g(\boldsymbol{\mu}, \mathbf{x})^T \right] = S^{-1} \underbrace{\mathbb{E}_{\mathbf{x}} \left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right]}_S S^{-1} = S^{-1}$$

- The Fisher kernel is then given as:

$$k(\mathbf{x}, \mathbf{x}') = g(\boldsymbol{\mu}, \mathbf{x})^T \mathbf{F}^{-1} g(\boldsymbol{\mu}, \mathbf{x}') = (\mathbf{x} - \boldsymbol{\mu})^T S^{-1} (\mathbf{x}' - \boldsymbol{\mu})$$

- This is the Mahalanobis distance!

Sigmoidal Kernel

- Another kernel function is the sigmoidal kernel

$$k(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^T \mathbf{x}' + b)$$

- Unfortunately the Gram matrix in general is not positive semidefinite.
- This kernel was used because it gives kernel expansions such as the SVM a superficial resemblance to neural network models.
- As we shall see later in this lecture, in the limit of an infinite number of basis functions, a Bayesian neural network with an appropriate prior reduces to a Gaussian process, thereby providing a link between neural networks and kernel methods.

- Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer



Radial Basis Functions

- By definition

$$\phi_j(x) = h(\|x - \mu_j\|)$$

where $\|x - \mu_j\|$ is a radial distance from a center μ_j .

- Originally were introduced for the problem of exact interpolation $f(x_n) = t_n, n = 1, \dots, N$

$$f(x) = \sum_{n=1}^N w_n h(\|x - x_n\|)$$

- There are as many unknowns w_n as target values t_n . Thus least squares leads to exact interpolation.
- If the target values t_n are noisy, exact interpolation overfits.

- Powell, M. J. D. (1987). [Radial basis functions for multivariable interpolation: a review](#). In J. C. Mason and M. G. Cox (Eds.), *Algorithms for Approximation*, pp. 143–167. Oxford University Press.



Radial Basis: Regularized Least Squares

- For a sum-of-squares error function with a regularizer defined in terms of a differential operator, the optimal solution is given by an **expansion in the Green's functions of the operator** (analogous to eigenvectors of a discrete matrix) again **with one basis function centered on each data point.**
- Green's functions for an isotropic differential operator as regularizer lead to radial basis functions (depending only on the distance from the data points)
- Interpolation is not exact due to the presence of the regularizer.

- Poggio, T. and F. Girosi (1990). [Networks for approximation and learning](#). *Proceedings of the IEEE* **78**(9), 1481–1497.
- Bishop, C. M. (1995a). [Neural Networks for Pattern Recognition](#). Oxford University Press



Interpolation with Noisy Inputs

- Consider the interpolation problem with noisy inputs rather than target values:

$$E = \frac{1}{2} \sum_{n=1}^N \int \left\{ y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n \right\}^2 \nu(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

- We can optimize with respect to the function $y(\mathbf{x})$ to obtain:

$$y(\mathbf{x}) = \sum_{n=1}^N t_n h(\mathbf{x} - \mathbf{x}_n)$$

- The basis functions are given as (Nadaraya-Watson model):

$$h(\mathbf{x} - \mathbf{x}_n) = \frac{\nu(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N \nu(\mathbf{x} - \mathbf{x}_n)}$$

- If noise distribution $\nu(\boldsymbol{\xi})$ is isotropic (function of $\|\boldsymbol{\xi}\|$), these basis functions would be radial.

- Bishop, C. M. (1995a). [Neural Networks for Pattern Recognition](#). Oxford University Press
- Webb, A. R. (1994). [Functional approximation by feed-forward networks: a least-squares approach to generalisation](#). *IEEE Transactions on Neural Networks* 5(3), 363–371.



Radial Basis Functions: Proof

- We take a variation of $y(\mathbf{x})$: $y(\mathbf{x}) \rightarrow y(\mathbf{x}) + \varepsilon\eta(\mathbf{x})$. Then:

$$E[y + \varepsilon\eta] = \frac{1}{2} \sum_{n=1}^N \int \left\{ y(\mathbf{x}_n + \boldsymbol{\xi}) + \varepsilon\eta(\mathbf{x}_n + \boldsymbol{\xi}) - t_n \right\}^2 v(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

- We expand in ε and set the coefficient in ε (functional 1st derivative) to zero:

$$\sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\} \eta(\mathbf{x}_n + \boldsymbol{\xi}) v(\boldsymbol{\xi}) d\boldsymbol{\xi} = 0$$

- Since this must be true for any variation $\eta(\mathbf{x})$, we choose:
 $\eta(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{z})$. This leads to:

$$\sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\} \delta(\mathbf{x}_n + \boldsymbol{\xi} - \mathbf{z}) v(\boldsymbol{\xi}) d\boldsymbol{\xi} = \sum_{n=1}^N \{y(\mathbf{z}) - t_n\} v(\mathbf{z} - \mathbf{x}_n) = 0 \Rightarrow$$

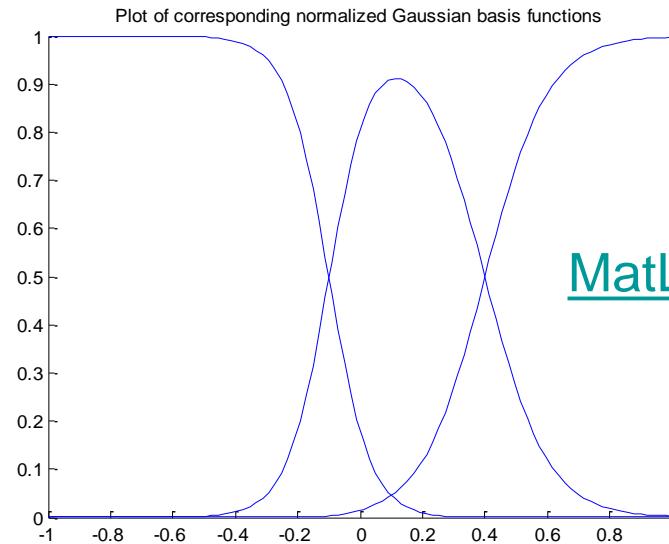
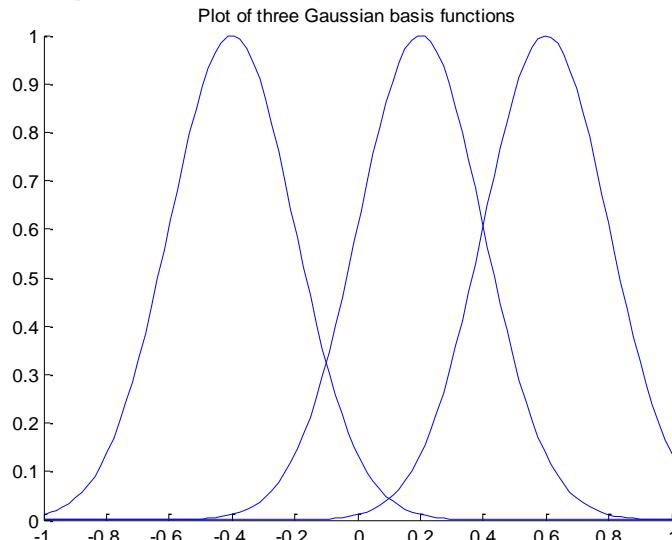
$$y(\mathbf{z}) = \sum_{n=1}^N t_n \left(v(\mathbf{z} - \mathbf{x}_n) / \sum_{l=1}^N v(\mathbf{z} - \mathbf{x}_l) \right)$$

Nadaraya-Watson Model

- The functions in the Nadaraya-Watson model are normalized:

$$h(x - x_n) = \frac{\nu(x - x_n)}{\sum_{n=1}^N \nu(x - x_n)}, \quad \sum_n h(x - x_n) = 1$$

- It is useful to have normalized basis functions to avoid regions in an input space where all of the basis functions take small values (leading to small predictions or predictions controlled solely by the bias parameter).



[MatLab Code](#)

Radial Basis For Regression

- In the application of kernel density estimation to regression, there is one basis function associated with every data point, and the corresponding model can be computationally costly to evaluate when making predictions for new data points.
- Models have therefore been proposed using the expansion in radial basis functions but where the number M of basis functions is smaller than the number N of data points.
- Typically, the number of basis functions, and the centers μ_i are determined based on the input data $\{x_n\}$ alone. The basis functions are kept fixed and the coefficients $\{w_i\}$ are determined by least squares.

- Broomhead, D. S. and D. Lowe (1988). [Multivariable functional interpolation and adaptive networks](#). *Complex Systems* **2**, 321–355.
- Moody, J. and C. J. Darken (1989). [Fast learning in networks of locally-tuned processing units](#). *Neural Computation* **1**(2), 281–294.
- Poggio, T. and F. Girosi (1990). [Networks for approximation and learning](#). *Proceedings of the IEEE* **78**(9), 1481–1497.



Reducing the Size of the Basis

- To reduce the algorithmic cost, we can keep the number of basis functions M smaller than input data size N
- Centers locations μ_i are determined based on the input data $\{x_n\}$ alone
- Coefficients $\{w_i\}$ are determined by least squares
- For the choice of centers:
 - Random subset of the training data points
 - Orthogonal least squares (sequential) – the next data point to be chosen as a basis function center gives the greatest error reduction
 - Using clustering algorithms (K-means) – basis functions centers no longer coincide with the training data points.
- Chen, S., C. F. N. Cowan, and P. M. Grant (1991). [Orthogonal least squares learning algorithm for radial basis function networks. *IEEE Transactions on Neural Networks* 2\(2\), 302–309.](#)



Parzen Density Estimator

- We have seen the prediction of the linear regression model - linear combination of t_n with 'equivalent kernel' values
- The same result is obtained from the Parzen density estimator.
 $f(\mathbf{x}, t)$ is the component density function centered at each training data point

$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n)$$

- Regression function

$$\begin{aligned} y(\mathbf{x}) &= \mathbb{E}[t | \mathbf{x}] = \int_{-\infty}^{+\infty} tp(t | \mathbf{x}) dt = \\ &= \frac{\int tp(\mathbf{x}, t) dt}{\int p(\mathbf{x}, t) dt} = \frac{\sum_n \int tf(\mathbf{x} - \mathbf{x}_n, t - t_n) dt}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt} \end{aligned}$$



Kernel Regression: Nadaraya-Watson Model

- Assume that the component density functions have zero mean so that

$$\int_{-\infty}^{+\infty} tf(x, t) dt = 0$$

for all values of x . Then by variable change

$$y(x) = \frac{\sum_n \int (t - t_n + t_n) f(x - x_n, t - t_n) dt}{\sum_m \int f(x - x_m, t - t_m) dt} = \frac{\sum_n g(x - x_n) t_n}{\sum_m g(x - x_m)} = \\ = \sum k(x, x_n) t_n$$

with $g(x) = \int_{-\infty}^{+\infty} f(x, t) dt$ and $k(x, x_n)$ given by

$$k(x, x_n) = \frac{g(x - x_n)}{\sum_m g(x - x_m)}$$



Kernel Regression: Nadaraya-Watson Model

- For localized kernels, this regression model gives more weight to the data points \mathbf{x}_n close to \mathbf{x} .

$$y(\mathbf{x}) = \sum_n k(\mathbf{x}, \mathbf{x}_n) t_n$$

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)}$$

- The kernel satisfies the summation constraint:

$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1$$

- In addition to the conditional expectation $y(\mathbf{x})$, the model provides the full conditional distribution:

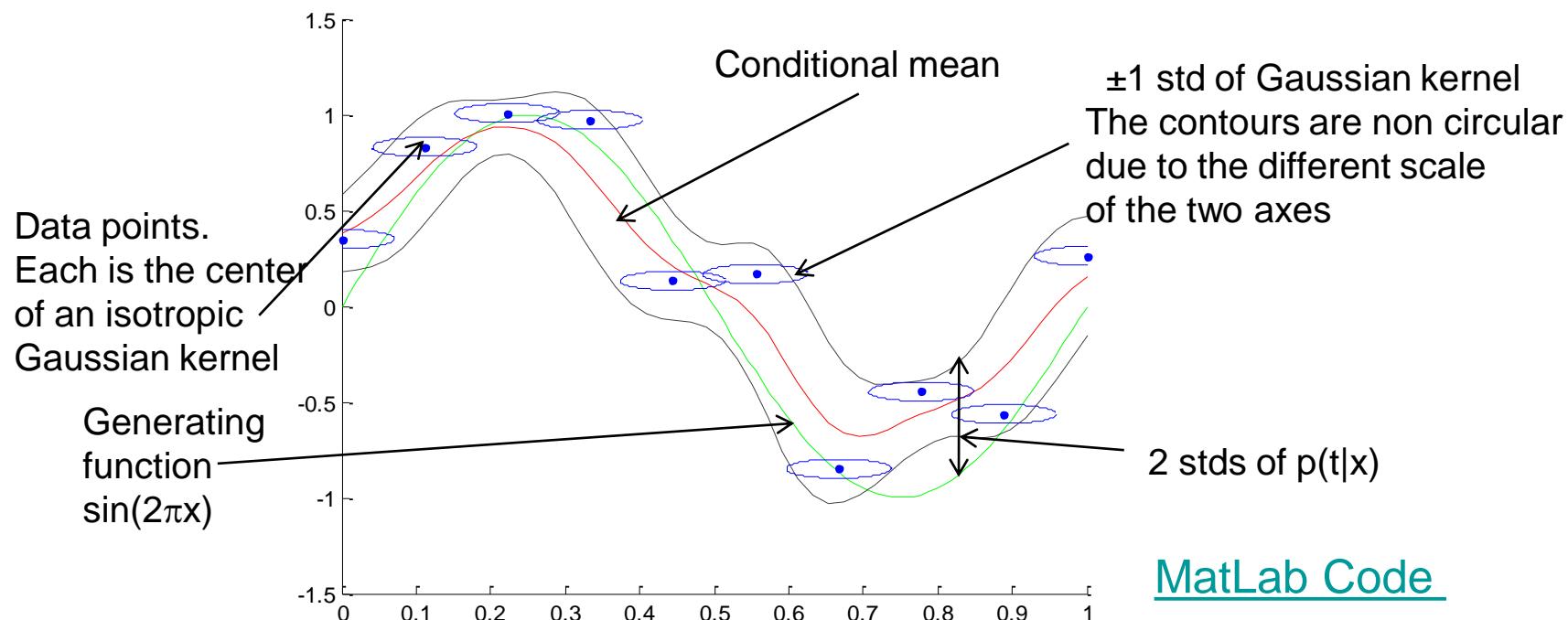
$$p(t | \mathbf{x}) = \frac{p(t, \mathbf{x})}{\int p(t, \mathbf{x}) dt} = \frac{\sum_n f(\mathbf{x} - \mathbf{x}_n, t - t_n)}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt}$$

- For Gaussian component distributions, this is a mixture of Gaussians.



Kernel Regression: Nadaraya-Watson Model

- Single input variable x . $f(x, t)$ is a zero-mean isotropic Gaussian kernels over the variable $z = (x, t)$ with variance σ^2



- Nadaraya, E. A. (1964). [On estimating regression. Theory of Probability and its Applications](#) **9**(1), 141–142.
- Watson, G. S. (1964). [Smooth regression analysis. Sankhya: The Indian Journal of Statistics. Series A](#) **26**, 359–372.

Kernel Regression: Nadaraya-Watson Model

- We can use more flexible forms of Gaussian components, e.g. having different variance parameters for the input and target variables.
- We could alternatively model $p(t, x)$ using a Gaussian mixture model, trained using EM techniques and then find $p(t|x)$.
 - ✓ In this case we don't have a representation in terms of kernel functions evaluated at the training set data points.
 - ✓ However, the number of components in the mixture model can be smaller than the number of training set points, resulting in a model that is faster to evaluate for test data points.

- [Ghahramani, Z. and M. I. Jordan \(1994\). Supervised learning from incomplete data via an EM approach.](#) In J. D. Cowan, G. T. Tesauro, and J. Alspector (Eds.), [Advances in Neural Information Processing Systems, Volume 6](#), pp. 120–127.



Gaussian Processes

- The idea is similar to linear regression with a fixed set of basis functions:

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

- However, here we forget about the parametric model $y(\mathbf{x}, \mathbf{w})$
- Instead we take an infinite number of basis functions given by a probability distribution over functions.
- This is not difficult to model since we only have to consider the values of the functions at training/test data points.
- Gaussian process models include:
 - Kriging
 - ARMA (autoregressive moving average)
 - Kalman filters
 - Radial basis function networks , etc.

Gaussian Processes

- Models equivalent to Gaussian processes have been widely studied in many different fields.
- In the geostatistics, Gaussian process regression is known as *kriging*.
- Similarly, ARMA (autoregressive moving average) models
- Kalman filters and radial basis function networks can all be viewed as forms of Gaussian process models.

- Cressie, N. (1993). [*Statistics for Spatial Data*](#). Wiley.
- MacKay, D. J. C. (1998). [*Introduction to Gaussian processes*](#). In C. M. Bishop (Ed.), *Neural Networks and Machine Learning*, pp. 133–166. Springer.
- Williams, C. K. I. (1999). [*Prediction with Gaussian processes: from linear regression to linear prediction and beyond*](#). In M. I. Jordan (Ed.), *Learning in Graphical Models*, pp. 599–621. MIT Press.
- MacKay (2003)
- Rasmussen, C. E. (1996). [*Evaluation of Gaussian Processes and Other Methods for Non-Linear Regression*](#). Ph. D. thesis, University of Toronto.
- Rasmussen, C. E. and C. K. I. Williams (2006). [*Gaussian Processes for Machine Learning*](#). MIT Press.
- Gibbs, 1997;



Linear Regression Revisited

- Consider a model defined by linear combination of M fixed basis functions:

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

- A Gaussian prior distribution over the weight vector \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

governed by the hyperparameter (precision) α induces then a probability distribution over functions $y(\mathbf{x})$.



Linear Regression Revisited

- Evaluating $y(\mathbf{x})$ for a set of training data points $\mathbf{x}_1, \dots, \mathbf{x}_N$, we have a joint distribution of $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$

$$\mathbf{y} = \Phi \mathbf{w}, \text{ with elements } y_n = y(\mathbf{x}_n) = \mathbf{w}^T \phi(\mathbf{x}_n)$$

where Φ is the design matrix with elements $\Phi_{nk} = \phi_k(\mathbf{x}_n)$

- $p(\mathbf{y})$ is Gaussian, and its mean and covariance can be shown to be

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = 0$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T \Rightarrow \text{cov}[\mathbf{y}] = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K}$$

where \mathbf{K} is the Gram matrix with elements defined in terms of the kernel function $k(\mathbf{x}, \mathbf{x}')$ as

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

- Up to now we only took data points + prior but no target values.



Gaussian Processes

- The model we have seen is one example of a Gaussian process
- Gaussian process is defined as a probability distribution over functions $y(x)$ such that the set of values of $y(x)$ evaluated at an arbitrary set of points x_1, \dots, x_N jointly have a Gaussian distribution (in 2D, we have a Gaussian Random Field).
- Key point: **the joint distribution is defined completely by second-order statistics (mean, covariance)**
- Since usually the mean is taken to be zero, we only need the covariance, i.e., the kernel-function:

$$\mathbb{E}[y(x_n) y(x_m)] = k(x_n, x_m)$$

- Instead of choosing (a limited number of) basis functions, we can directly choose a kernel function (which may result in an infinite number of basis functions)

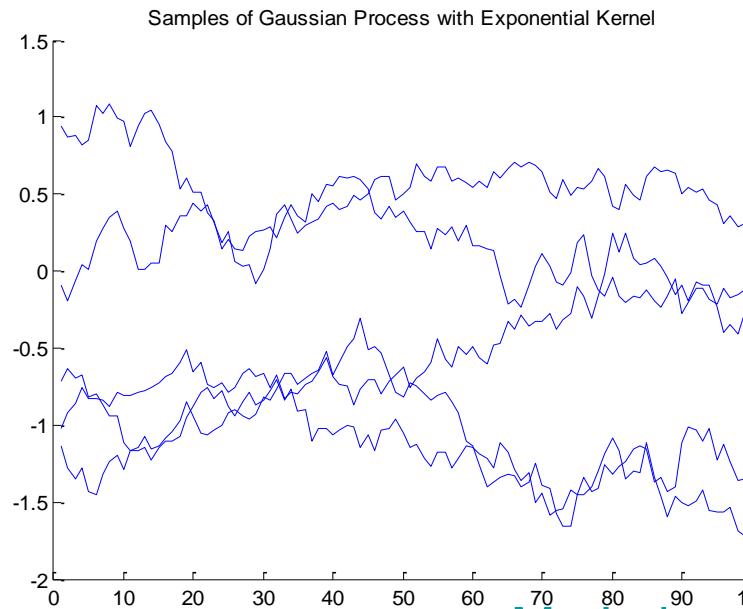
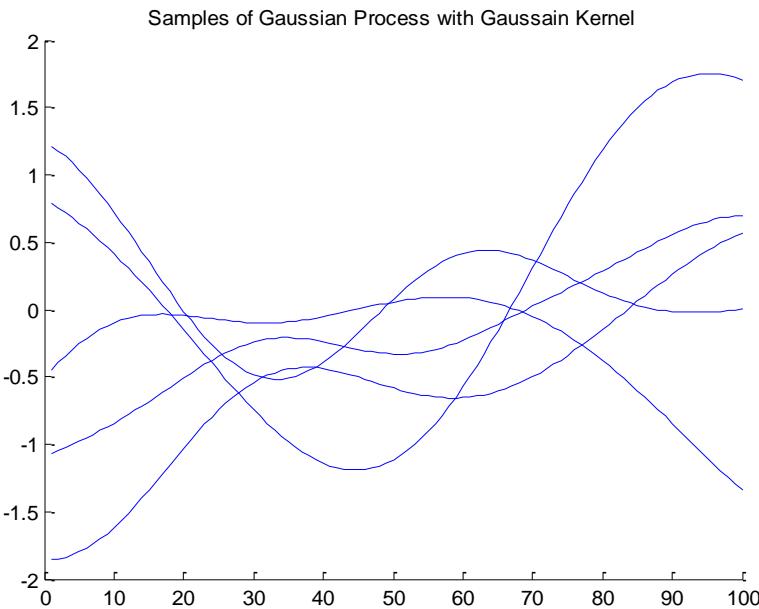


Gaussian Processes

- We show examples of functions drawn from Gaussian processes for a Gaussian kernel and an exponential kernel.

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

$k(\mathbf{x}, \mathbf{x}') = \exp(-\theta |\mathbf{x} - \mathbf{x}'|)$
(Ornstein-Uhlenbeck process)



MatLab code

- Uhlenbeck, G. E. and L. S. Ornstein (1930). [On the theory of Brownian motion](#). *Phys. Rev.* **36**, 823–841.



Gaussian Process Regression

- To use Gaussian processes for regression, we need to model noise:

$$t_n = y_n + \varepsilon_n, \quad \text{with} \quad y_n = y(\mathbf{x}_n)$$

- For noise processes with a Gaussian distribution we obtain

$$p(t_n | y_n) = \mathcal{N}(t_n | y_n, \beta^{-1})$$

- The joint distribution for $\mathbf{t} = (t_1, t_2, \dots, t_n)^T$ conditioned on $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ is then (since the noise is assumed to be independent for each data point)

$$p(\mathbf{t} | \mathbf{y}) = \mathcal{N}(\mathbf{t} | \mathbf{y}, \beta^{-1} \mathbf{I}_N)$$

Gaussian Process Regression

- From the definition of a Gaussian process we have

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})$$

- The kernel function that determines \mathbf{K} is chosen such that for points $\mathbf{x}_n, \mathbf{x}_m$ that are similar the corresponding values $y(\mathbf{x}_n), y(\mathbf{x}_m)$ will be stronger correlated.
- For the marginal distribution $p(\mathbf{t})$, we need to integrate over \mathbf{y} (see Appendix and earlier lecture notes):

$$p(\mathbf{t}) = \int p(\mathbf{t} | \mathbf{y}) p(\mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{C})$$

with $\mathbf{C} = \mathbf{K} + \beta^{-1} \mathbf{I}$ or $C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \delta_{nm}$

Appendix: Linear Gaussian Model

- Consider a linear Gaussian model: A Gaussian marginal distribution $p(\mathbf{x})$ and a Gaussian conditional distribution $p(\mathbf{y}|\mathbf{x})$ in which $p(\mathbf{y}|\mathbf{x})$ has a mean that is a linear function of \mathbf{x} , and a covariance which is independent of \mathbf{x} .

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

- We want to find $p(\mathbf{y})$. We start with the joint distribution over $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ which is a Gaussian and then integrate \mathbf{x} out.

- Finally we obtained:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

Practical Kernel Function

- One widely used kernel function for Gaussian process regression is given here:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

- The term involving θ_3 corresponds to a parametric model that is a linear function of the input variables.

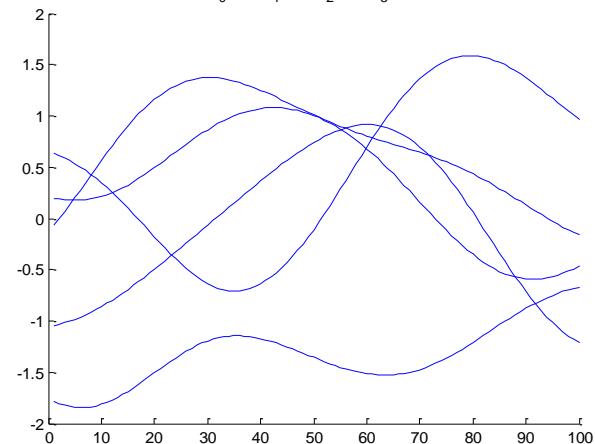


Distribution over Functions

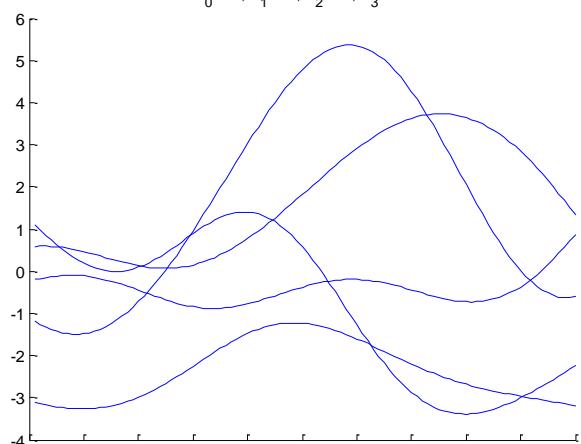
- Sample functions from prior $p(y)$ with common kernel function

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

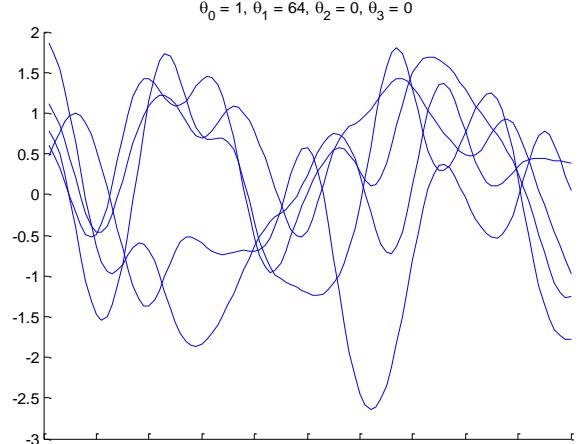
$$\theta_0 = 1, \theta_1 = 4, \theta_2 = 0, \theta_3 = 0$$



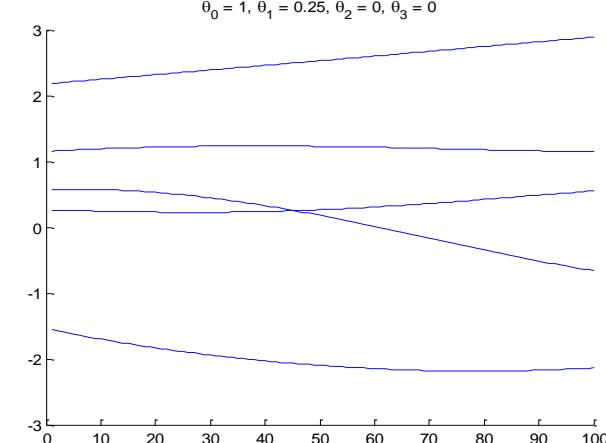
$$\theta_0 = 9, \theta_1 = 4, \theta_2 = 0, \theta_3 = 0$$



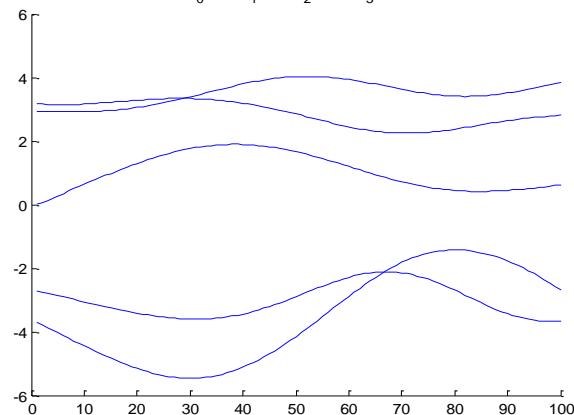
$$\theta_0 = 1, \theta_1 = 64, \theta_2 = 0, \theta_3 = 0$$



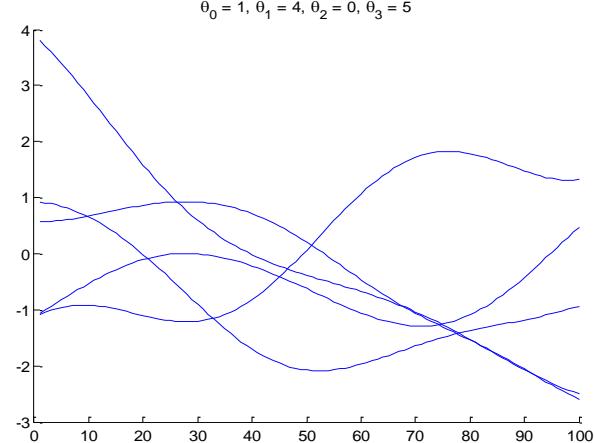
$$\theta_0 = 1, \theta_1 = 0.25, \theta_2 = 0, \theta_3 = 0$$



$$\theta_0 = 1, \theta_1 = 4, \theta_2 = 10, \theta_3 = 0$$



$$\theta_0 = 1, \theta_1 = 4, \theta_2 = 0, \theta_3 = 5$$



MatLab Code



Distribution Over Functions

- Sampling of data points $\{t_n\}$ from a Gaussian process.

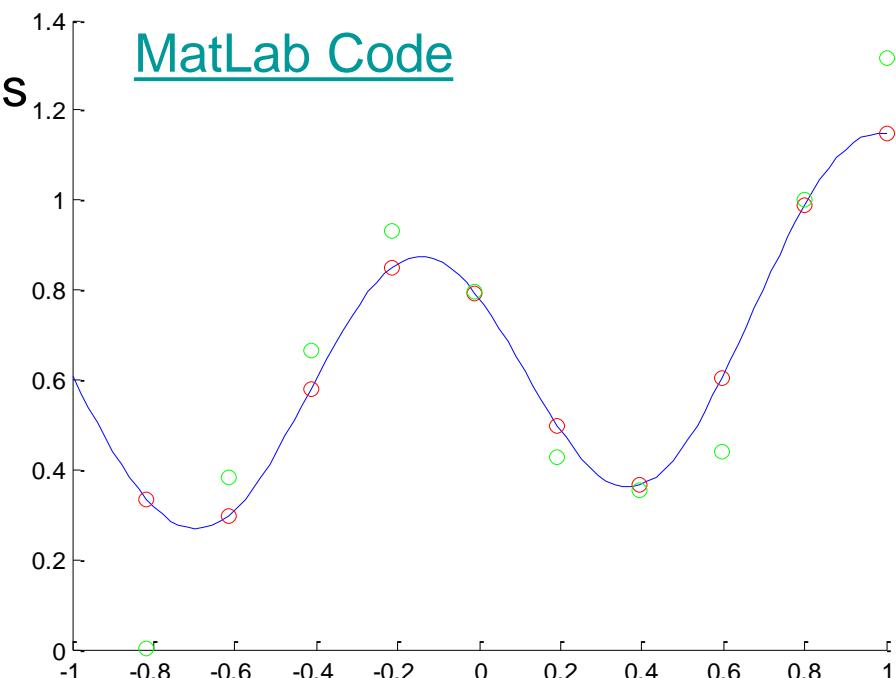
- The blue curve shows a sample function from the Gaussian process prior over functions.

$$p(y) = \mathcal{N}(y | \theta, K)$$

- The red points show the values of y_n by evaluating the function at a set of input values $\{x_n\}$.

- The corresponding values of $\{t_n\}$, shown in green, are obtained by adding independent Gaussian noise to each of the $\{y_n\}$.

$$p(t) = \mathcal{N}(t | 0, C)$$



Making Predictions

- So far we have a model for the joint probability distribution over sets of data points.
- For predictions t_{N+1} for a new input variable x_{N+1} , we need to evaluate the predictive distribution $p(t_{N+1}|\mathbf{t}_N)$.
- By partitioning the joint Gaussian distribution over x_1, \dots, x_N, x_{N+1} , we obtain $p(t_{N+1}|\mathbf{t})$ given by its mean and covariance (see following Appendix with results for linear Gaussian models discussed in an earlier lecture).
- The joint distribution over t_1, \dots, t_N, t_{N+1} is:

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

- We partition the covariance matrix as follows:



Making Predictions

$$p(\boldsymbol{t}_{N+1}) = \mathcal{N}(\boldsymbol{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}$$

- Here \mathbf{C}_N is the $N \times N$ covariance matrix with elements given by

$$\mathbf{C}_N(\boldsymbol{x}_n, \boldsymbol{x}_m) = k(\boldsymbol{x}_n, \boldsymbol{x}_m) + \beta^{-1} \delta_{nm}, n, m = 1, \dots, N$$

- The vector \mathbf{k} has elements $k(\boldsymbol{x}_n, \boldsymbol{x}_{N+1})$ for $n = 1, \dots, N$, and the scalar

$$c = k(\boldsymbol{x}_{N+1}, \boldsymbol{x}_{N+1}) + \beta^{-1}$$



Appendix: Conditional Gaussian Distributions

- If two sets of variables are jointly Gaussian, then the conditional distribution of one set conditioned on the other is again Gaussian.
- Suppose x is a D -dimensional vector with Gaussian distribution $\mathcal{N}(x|\mu, \Sigma)$ and that we partition x into two disjoint subsets x_a (M components) and x_b ($D - M$ components).

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$



Appendix: Conditional Gaussian Distributions

- This partition also implies similar partitions for the mean and covariance.

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

- $\boldsymbol{\Sigma}^T = \boldsymbol{\Sigma}$ implies that $\boldsymbol{\Sigma}_{aa}$ and $\boldsymbol{\Sigma}_{bb}$ are symmetric and

$$\boldsymbol{\Sigma}_{ba} = \boldsymbol{\Sigma}_{ab}^T$$



Appendix: The Precision Matrix

- We define the precision matrix Λ as Σ^{-1} .
- Its partition is given as above

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

where from $\Sigma^T = \Sigma$ we conclude that Λ_{aa} and Λ_{bb} are symmetric and

$$\Lambda_{ba} = \Lambda_{ab}^T$$

- Note that the above partition does NOT imply that Λ_{aa} is the inverse of Σ_{aa} , etc.



Appendix: The Conditional Distribution

- We are now interested to compute $p(\mathbf{x}_a|\mathbf{x}_b)$. An easy way to do that is to look at $p(\mathbf{x}_a, \mathbf{x}_b)$ considering \mathbf{x}_b constant.
- We have proved that the conditional is Gaussian with mean and variance given as:

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}$$

- Note that the conditional mean is linear in \mathbf{x}_b and the conditional variance is independent of \mathbf{x}_b .

Making Predictions

- Using the results from the Appendix, $p(t_{N+1}|\mathbf{t})$ is a Gaussian with mean and covariance:

$$\mathbf{t} = (t_1, \dots, t_N, t_{N+1})^T = (b \quad a)^T \quad \mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix} \quad p(t_{N+1}) = \mathcal{N}(t_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$
$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b) \Rightarrow \boldsymbol{\mu}_{t_{N+1}|\mathbf{t}_N} = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \Rightarrow \Sigma_{t_{N+1}|\mathbf{t}_N} = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}$$

- Finally we obtain $p(t_{N+1}|\mathbf{t})$ with both the mean and variance functions of \mathbf{x}_{N+1}

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}_{N+1}), \quad a_n = (\mathbf{C}_N^{-1} \mathbf{t}_N)_n,$$

$$\mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}, \quad c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

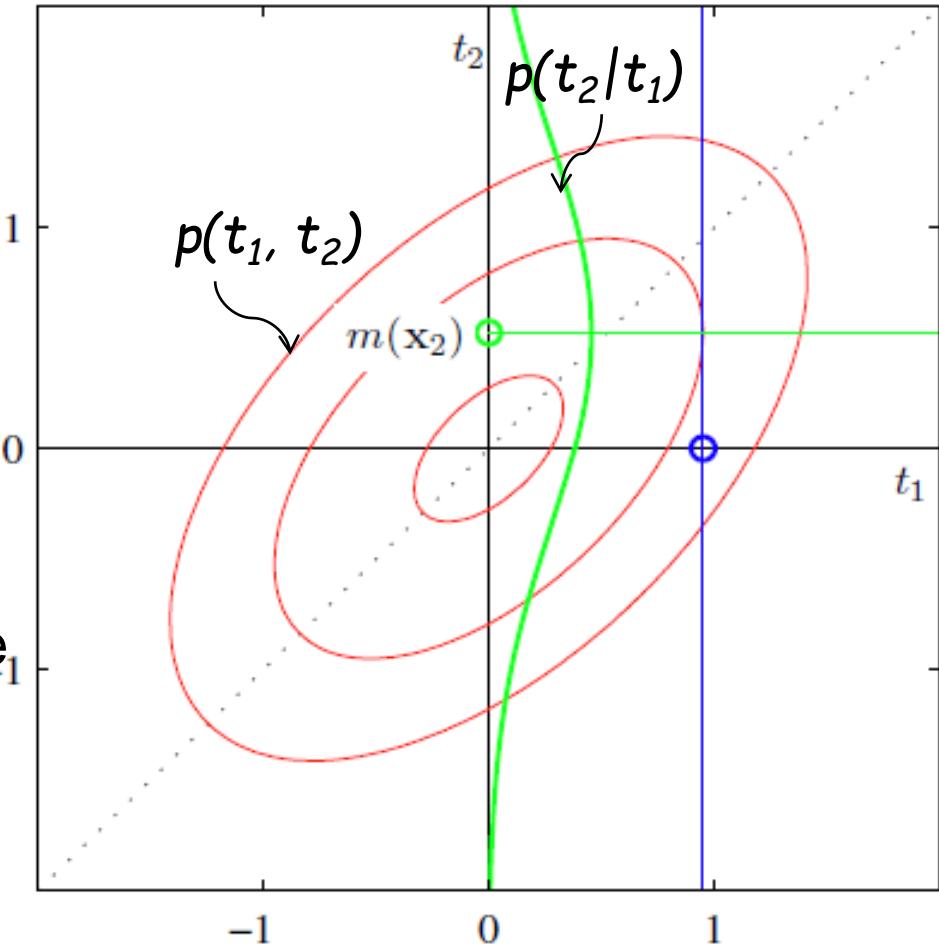
- If $k(\mathbf{x}_n, \mathbf{x}_m)$ depends only on $\|\mathbf{x}_n - \mathbf{x}_m\|$, then we obtain an expansion in radial basis functions.



Gaussian Process Regression

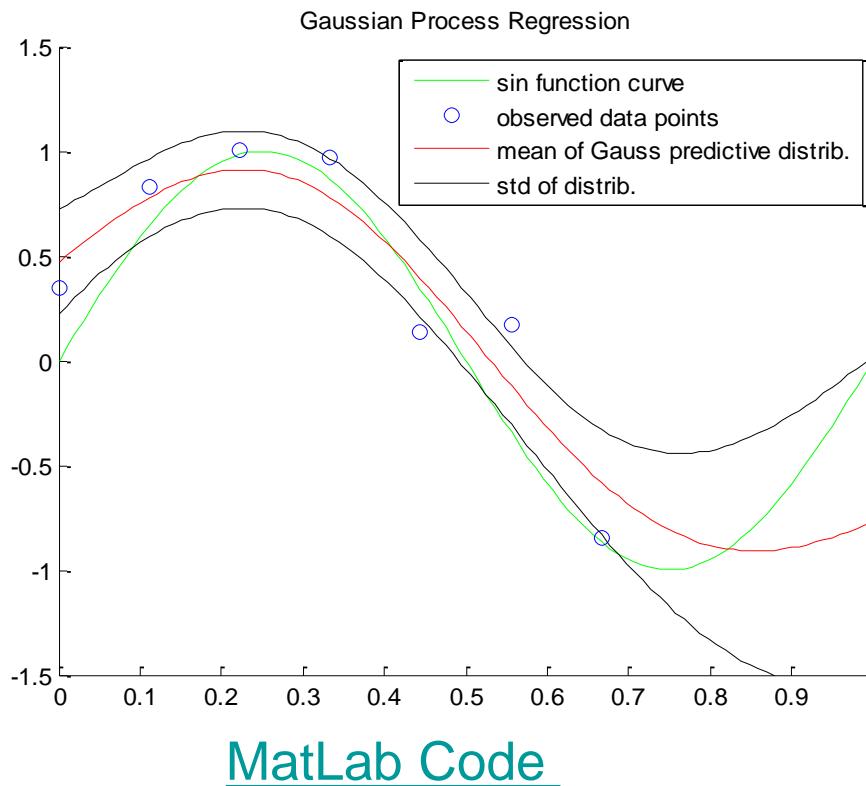
- Illustration of the mechanism of Gaussian process regression for the case of one training point and one test point.

- The red ellipses show contours of the joint distribution $p(t_1, t_2)$.
- t_1 is the training data point, and conditioning on the value of t_1 (blue line), we obtain $p(t_2|t_1)$ shown as a function of t_2 by the green curve.



Making Predictions

- Green curve: original sinusoidal function; blue points: sampled training data points with additional noise; red line: mean estimate; black lines: $\pm 2\sigma$.
- Note the increase of uncertainty in the right of the data points.



Restrictions on the kernel Matrix

- The covariance matrix must be positive definite.

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \delta_{nm}, n, m = 1, \dots, N$$

- If λ_i is an eigenvalue of K , then the corresponding eigenvalue of C will be $\lambda_i + \beta^{-1}$.
- It is therefore sufficient that the kernel matrix $k(\mathbf{x}_n, \mathbf{x}_m)$ be positive semidefinite for any pair of points \mathbf{x}_n and \mathbf{x}_m , so that $\lambda_i > 0$. Any $\lambda_i = 0$ will still give rise to a positive eigenvalue for C because $\beta > 0$.
- This is the same restriction on the kernel function [discussed earlier](#), and so we can exploit the techniques discussed there.

Making Predictions

- Note that the mean of the conditional $p(t_{N+1} | \mathbf{t})$

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N, \quad \mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T$$

can be written as:

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}_{N+1}),$$

where

$$\mathbf{C}_N^{-1} \mathbf{t}_N = (a_1 \dots a_N)^T$$

- Thus if the kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$ depends only on the distance $\|\mathbf{x}_n - \mathbf{x}_m\|$ we obtain an expansion in radial basis.



Making Predictions

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N, \quad \mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T$$
$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}, \quad c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

- These results define the predictive distribution for Gaussian process regression with an arbitrary kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$
- If the kernel function $k(\mathbf{x}, \mathbf{x}')$ is defined in terms of a finite set of basis functions, the results become identical to those we discussed in an earlier lectures on regression (see proof in the next two slides)

$$p(t | x, \mathbf{x}, \mathbf{t}) = \int p(t | x, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}\left(t | \mathbf{m}_N^T \boldsymbol{\phi}(x), \sigma_N^2(x)\right)$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

$$\sigma_N^2(x) = \frac{1}{\beta} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)$$

Kernels Defined by a Finite Basis

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N, \quad \mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T$$

□ We can write the vector \mathbf{k} as:

$$\mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T = \frac{1}{\alpha} \begin{pmatrix} \boldsymbol{\phi}^T(\mathbf{x}_1) \\ \boldsymbol{\phi}^T(\mathbf{x}_2) \\ \vdots \\ \boldsymbol{\phi}^T(\mathbf{x}_N) \end{pmatrix} \boldsymbol{\phi}(\mathbf{x}_{N+1}) = \frac{1}{\alpha} \boldsymbol{\Phi} \boldsymbol{\phi}(\mathbf{x}_{N+1})$$

and use

$$\mathbf{C}_N = \alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T + \beta^{-1} \mathbf{I}_N$$

□ Using this, the predicted mean takes the form

$$\begin{aligned} m(\mathbf{x}_{N+1}) &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N = \frac{1}{\alpha} \boldsymbol{\phi}^T(\mathbf{x}_{N+1}) \boldsymbol{\Phi}^T (\alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T + \beta^{-1} \mathbf{I}_N)^{-1} \mathbf{t}_N = \\ &\alpha^{-1} \boldsymbol{\phi}^T(\mathbf{x}_{N+1}) \beta \boldsymbol{\Phi}^T \left(\begin{array}{cc} \alpha^{-1} \boldsymbol{\Phi} & \beta \boldsymbol{\Phi}^T + \mathbf{I}_N \\ \mathbf{B} & \mathbf{A} \end{array} \right)^{-1} \mathbf{t}_N \quad (I + AB)^{-1} A = A(I + BA)^{-1} = \\ &\alpha^{-1} \boldsymbol{\phi}^T(\mathbf{x}_{N+1}) (\beta \boldsymbol{\Phi}^T \alpha^{-1} \boldsymbol{\Phi} + \mathbf{I}_M)^{-1} \beta \boldsymbol{\Phi}^T \mathbf{t}_N = \\ &\boldsymbol{\phi}^T(\mathbf{x}_{N+1}) \beta (\beta \boldsymbol{\Phi} \boldsymbol{\Phi}^T + \alpha \mathbf{I}_M)^{-1} \boldsymbol{\Phi}^T \mathbf{t}_N = \beta \boldsymbol{\phi}^T(\mathbf{x}_{N+1}) \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}_N, \text{ the required result} \end{aligned}$$

Kernels Defined by a Finite Basis

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}, \quad c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

- Using $\mathbf{k} = \alpha^{-1} \Phi \phi(\mathbf{x}_{N+1})$, $c = \alpha^{-1} \phi(\mathbf{x}_{N+1})^T \phi(\mathbf{x}_{N+1}) + \beta^{-1}$ and $\mathbf{C}_N = \alpha^{-1} \Phi \Phi^T + \beta^{-1} \mathbf{I}_N$, we compute the predictive variance as:

$$\begin{aligned}\sigma^2(\mathbf{x}_{N+1}) &= \alpha^{-1} \phi(\mathbf{x}_{N+1})^T \phi(\mathbf{x}_{N+1}) + \beta^{-1} - \phi(\mathbf{x}_{N+1})^T \Phi^T \alpha^{-1} \left(\alpha^{-1} \Phi \Phi^T + \beta^{-1} \mathbf{I}_N \right)^{-1} \alpha^{-1} \Phi \phi(\mathbf{x}_{N+1}) = \\ \beta^{-1} + \phi(\mathbf{x}_{N+1})^T \left\{ \underbrace{\alpha^{-1} \mathbf{I}_M}_{\mathbf{A}^{-1}} - \underbrace{\Phi^T}_{\mathbf{B}} \underbrace{\alpha^{-1} \mathbf{I}_M}_{\mathbf{A}^{-1}} \left(\underbrace{\alpha^{-1} \Phi \Phi^T}_{\mathbf{C}\mathbf{A}^{-1}\mathbf{B}} + \underbrace{\beta^{-1} \mathbf{I}_N}_{\mathbf{D}} \right)^{-1} \underbrace{\Phi}_{\mathbf{C}} \underbrace{\alpha^{-1} \mathbf{I}_M}_{\mathbf{A}^{-1}} \right\} \phi(\mathbf{x}_{N+1}) = \\ \beta^{-1} + \phi(\mathbf{x}_{N+1})^T (\alpha \mathbf{I}_M + \beta \Phi^T \Phi)^{-1} \phi(\mathbf{x}_{N+1}) &= \beta^{-1} + \phi(\mathbf{x}_{N+1})^T S_N \phi(\mathbf{x}_{N+1}), \text{ the required result}\end{aligned}$$

- Here we used the Woodbury identity with \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} as shown.

$$(\mathbf{A} + \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{D} + \mathbf{C} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{C} \mathbf{A}^{-1}$$

Gaussian process Vs. Basis Function Model

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N, \quad \sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}, \quad c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

$$p(t | x, \mathbf{x}, \mathbf{t}) = \mathcal{N}\left(t | \mathbf{m}_N^T \boldsymbol{\phi}(x), \sigma_N^2(x)\right) \quad \begin{aligned} \mathbf{m}_N &= \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi} \end{aligned} \quad \sigma_N^2(x) = \frac{1}{\beta} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)$$

- In Gaussian processes, we invert \mathbf{C}_N of size $N \times N$, cost $\mathcal{O}(N^3)$
- In the basis function model, we invert \mathbf{S}_N , size $M \times M$, cost $\mathcal{O}(M^3)$
- These matrix inversions must be performed once for the given training set.
- For each new test point, both methods require a vector-matrix multiplications ($\mathcal{O}(N^2)$ vs. $\mathcal{O}(M^2)$). **If $M < N$, it will be computationally more efficient to work in the basis function framework.**
- An advantage of Gaussian processes is that we can consider covariance functions that can only be expressed in terms of an infinite number of basis functions.



Large Training Data Sets

- For large training data sets, GPs can become infeasible, and approximation schemes have been developed with better scaling with training set size than the exact approach.

- Gibbs, M. N. (1997). *Bayesian Gaussian processes for regression and classification*. Phd thesis, University of Cambridge.
- Tresp, V. (2001). *Scaling kernel-based systems to large data sets*. *Data Mining and Knowledge Discovery* 5(3), 197–211.
- Smola, A. J. and P. Bartlett (2001). *Sparse greedy Gaussian process regression*. In T. K. Leen, T. G. Dietterich, and V. Tresp (Eds.), *Advances in Neural Information Processing Systems*, Volume 13, pp. 619–625. MIT Press.
- Williams, C. K. I. and M. Seeger (2001). *Using the Nyström method to speed up kernel machines*. In T. K. Leen, T. G. Dietterich, and V. Tresp (Eds.), *Advances in Neural Information Processing Systems*, Volume 13, pp. 682–688. MIT Press.
- Csató, L. and M. Opper (2002). *Sparse on-line Gaussian processes*. *Neural Computation* 14(3), 641– 668.
- Seeger, M., C. K. I. Williams, and N. Lawrence (2003). *Fast forward selection to speed up sparse Gaussian processes*. In C. M. Bishop and B. Frey (Eds.), *Proceedings Ninth International Workshop on Artificial Intelligence and Statistics*, Key West, Florida.
- Bishop, C. M. and I. T. Nabney (2008). *Pattern Recognition and Machine Learning: A Matlab Companion*. Springer. In preparation.



Extension of Gaussian Processes

- Extension of GPs to multiple target variables (co-kriging) is straightforward.
- Other extensions of GPs include
 - modelling the distribution over low-dimensional manifolds for unsupervised learning and
 - the solution of SPDEs

- Cressie, N. (1993). *Statistics for Spatial Data*. Wiley.
- Bishop, C. M., M. Svensén, and C. K. I. Williams (1998a). *Developments of the Generative Topographic Mapping*. *Neurocomputing* 21, 203– 224.
- Graepel, T. (2003). *Solving noisy linear operator equations by Gaussian processes: Application to ordinary and partial differential equations*. In *Proceedings of the Twentieth International Conference on Machine Learning*, pp. 234–241.



Learning Hyperparameters θ

- Instead of fixing the covariance function, use a parametric family of functions and infer the parameters from the data.
- Parameters θ : length scale of correlations, and precision of noise β .
- Techniques for learning θ are based on the likelihood $p(\mathbf{t}|\theta)$ – analogous to model evidence techniques used earlier.
- Simplest approach:
 - Maximizing the log-likelihood (MLE): $\ln p(\mathbf{t}|\theta)$
$$\ln p(\mathbf{t}|\theta) = -\frac{1}{2} \ln |\mathbf{C}_N| - \frac{1}{2} \mathbf{t}^T \mathbf{C}_N^{-1} \mathbf{t} - \frac{N}{2} \ln(2\pi)$$
 - Use gradient optimization
$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\theta) = -\frac{1}{2} \text{tr} \left[\mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \right] + \frac{1}{2} \mathbf{t}^T \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \mathbf{C}_N^{-1} \mathbf{t}$$
 - $\ln p(\mathbf{t}|\theta)$ is in general non-convex and can have multiple maxima

- Fletcher, R. (1987). *Practical Methods of Optimization* (Second ed.). Wiley.
- Nocedal, J. and S. J. Wright (1999). *Numerical Optimization*. Springer.



Appendix

□ Show that

$$\frac{\partial \mathbf{A}^{-1}}{\partial x} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$

Indeed note that

$$\frac{\partial}{\partial x} (\mathbf{A}\mathbf{A}^{-1}) = \mathbf{0} \Rightarrow \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} + \mathbf{A} \frac{\partial \mathbf{A}^{-1}}{\partial x} = \mathbf{0} \Rightarrow \frac{\partial \mathbf{A}^{-1}}{\partial x} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$

□ Show that (\mathbf{A} symmetric)

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$

Indeed using the eigen-decomposition of \mathbf{A} , note that

$$\mathbf{A} = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^T, \quad \mathbf{A}^{-1} = \sum_i \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T, \quad \text{tr} \mathbf{A} = \sum_i \lambda_i, \quad |\mathbf{A}| = \prod_i \lambda_i \Rightarrow$$

$$\ln |\mathbf{A}| = \sum_i \ln \lambda_i \Rightarrow \frac{\partial}{\partial x} \ln |\mathbf{A}| = \sum_i \frac{\partial \ln \lambda_i}{\partial x} = \sum_i \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x} = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$



Learning Hyperparameters θ

- One can also introduce a prior $p(\theta)$ and maximize the log-posterior (maximum a posteriori): $\ln p(t|\theta) + \ln p(\theta)$
- In a Bayesian setting, we need to compute marginals over θ weighted with $p(t|\theta)p(\theta)$; this is not tractable → use approximations
- We have assumed that the contribution to the predictive variance arising from the noise β is constant. The noise might not be additive but dependent on x (**heteroscedastic**)

$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{nm}, n, m = 1, \dots, N$$

- A second Gaussian process can be introduced to represent the dependency of β on x .
 - Use a Gaussian process for $\ln \beta$ since $\beta \geq 0$.
 - Goldberg, P. W., C. K. I. Williams, and C. M. Bishop (1998). [Regression with input-dependent noise: A Gaussian process treatment](#). In [Advances in Neural Information Processing Systems](#), Volume 10, pp. 493–499. MIT Press.



Automatic Relevance Detection

- An additional hyperparameter can be introduced for each input dimension, e.g. in two-dimensions:

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp\left(-\frac{1}{2} \sum_{i=1}^2 \eta_i (x_i - x'_i)^2\right)$$

- Hyperparameter optimization by maximum likelihood allows then a different weighting of each dimension.
- Unrelevant dimensions (with small weights) can be detected and discarded.

Automatic Relevance Detection

$$k(x, x') = \theta_0 \exp\left(-\frac{1}{2} \sum_{i=1}^2 \eta_i (x_i - x'_i)^2\right)$$

- Samples from the resulting prior over functions $y(x)$ are shown next for two settings of η_i .
- As a particular η_i becomes small, the function becomes insensitive to x_i .
- By adapting these parameters to a data set using MLE, it becomes possible to detect input variables that have little effect on the predictive distribution (such inputs are discarded)

- MacKay, D. J. C. (1994). [Bayesian methods for backprop networks](#). In E. Domany, J. L. van Hemmen, and K. Schulten (Eds.), *Models of Neural Networks, III*, Chapter 6, pp. 211–254. Springer.
- Neal, R. M. (1996). [Bayesian Learning for Neural Networks](#). Springer. Lecture Notes in Statistics 118.

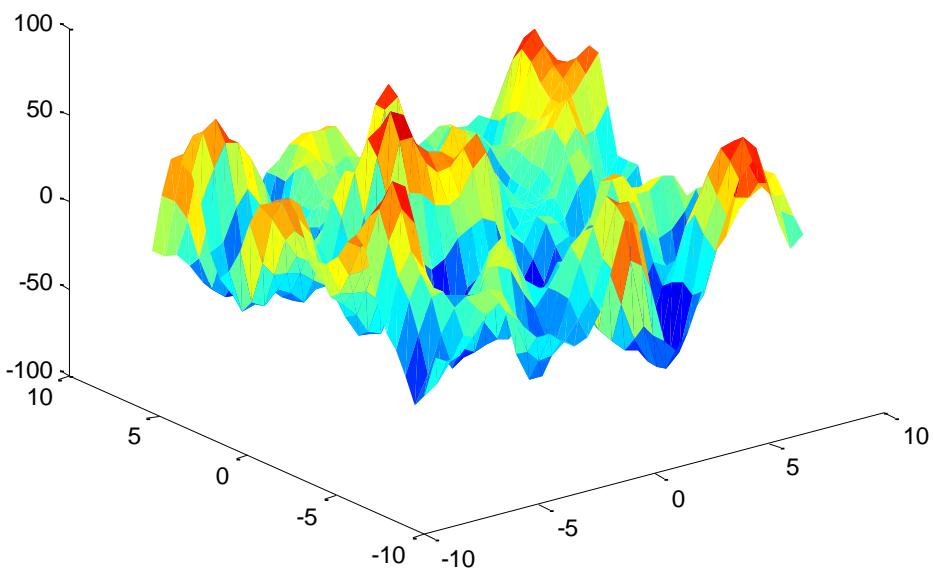


Automatic Relevance Detection

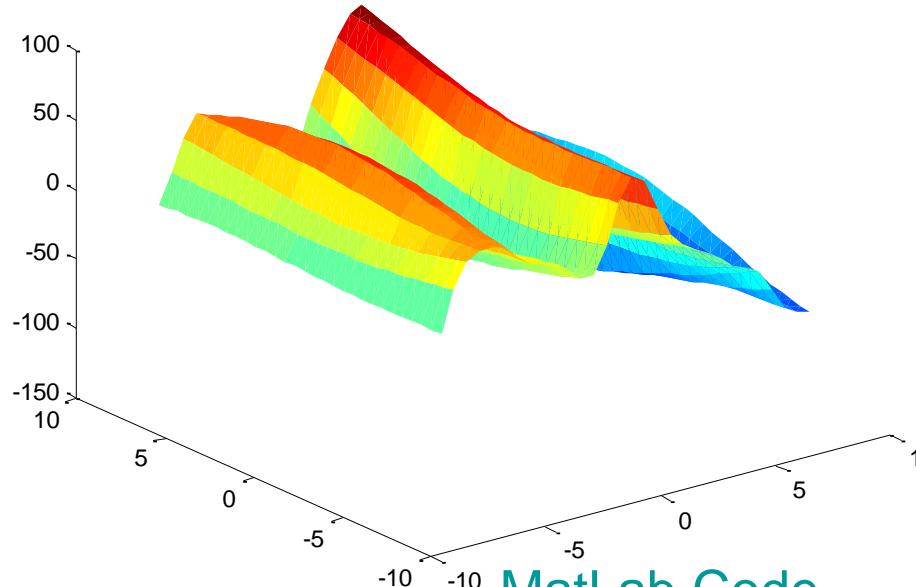
$$k(x, x') = \theta_0 \exp\left(-\frac{1}{2} \sum_{i=1}^2 \eta_i (x_i - x'_i)^2\right)$$

- Samples from the resulting prior over functions $y(x)$ are shown for two settings of η_i .

$$\eta_1 = \eta_2 = 1$$



$$\eta_1 = 1, \eta_2 = 0.01$$



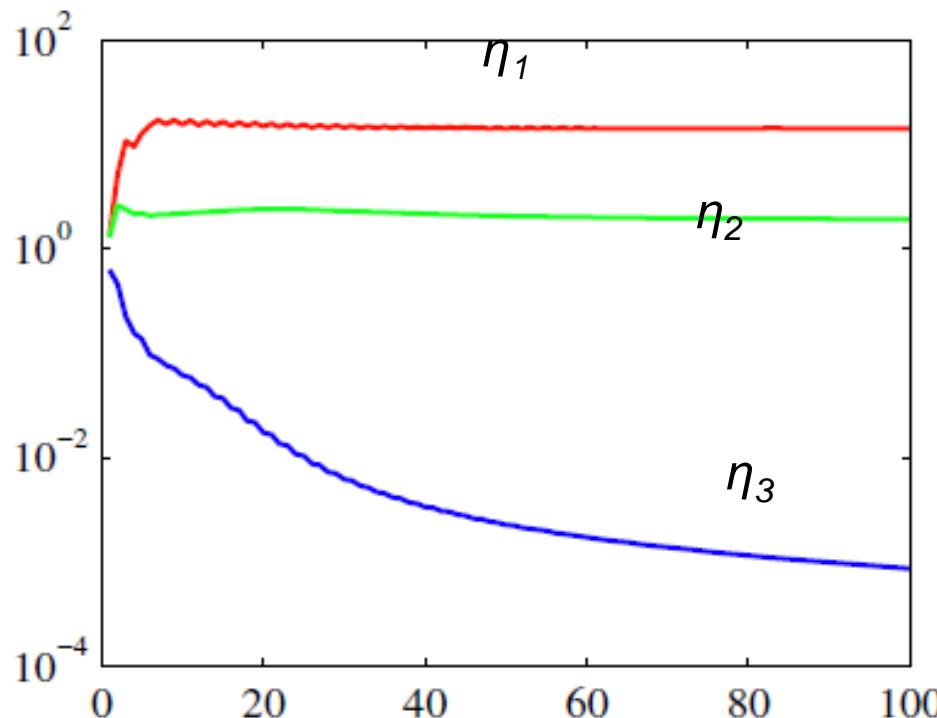
[MatLab Code](#)

Automatic Relevance Detection

- ARD is illustrated next using a data set having 3 inputs x_1, x_2 and x_3 .
- The target variable t , is generated by sampling 100 values of x_1 from a Gaussian, evaluating the function $\sin(2\pi x_1)$, and then adding Gaussian noise.
- Values of x_2 are given by copying the corresponding values of x_1 and adding noise, and values of x_3 are sampled from an independent Gaussian distribution.
- Thus x_1 is a good predictor of t , x_2 is a more noisy predictor of t , and x_3 has only chance correlating with t .
- The marginal likelihood for a Gaussian process with ARD parameters η_1, η_2, η_3 is optimized using **the scaled conjugate gradients algorithm**.

Automatic Relevance Detection

- We see (log scale in the vertical axis) that η_1 converges to a relatively large value, η_2 converges to a much smaller value, and η_3 becomes very small indicating that x_3 is irrelevant for predicting t .



[MatLab Code](#)

- Nabney, I. T. (2002). [*Netlab: Algorithms for Pattern Recognition*](#). Springer.



Automatic Relevance Detection

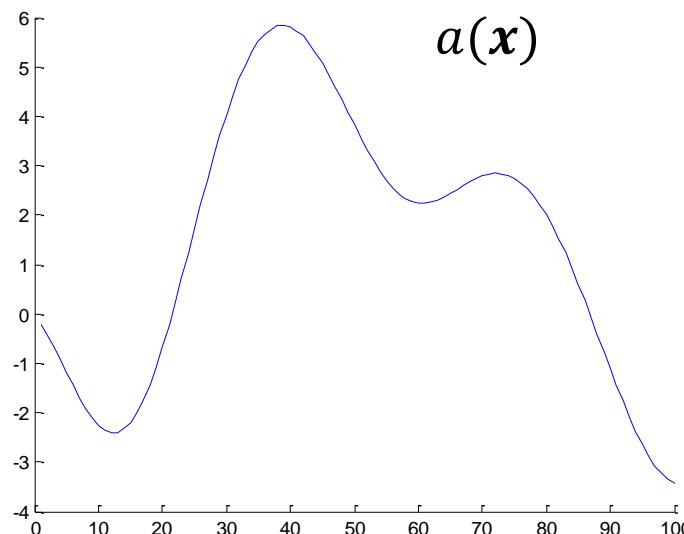
- The ARD framework can be applied to general kernels of the form:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \sum_{i=1}^D \eta_i (x_{ni} - x_{mi})^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

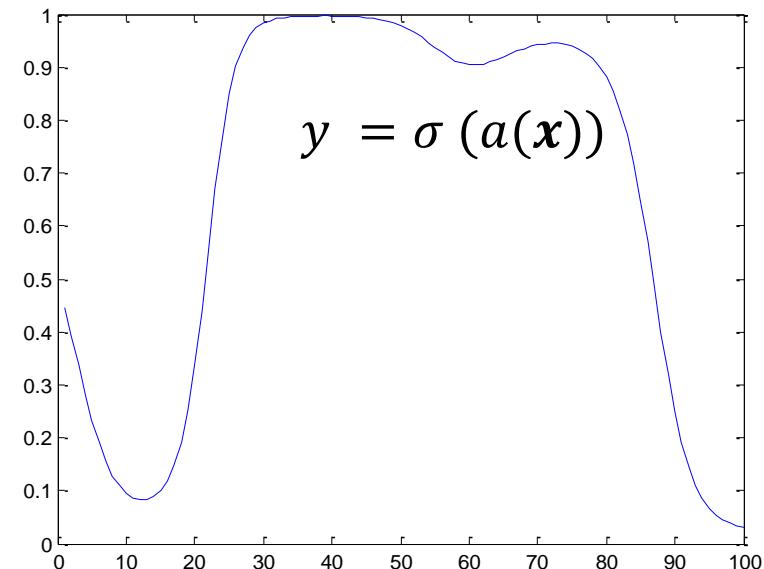
where D is the input dimensionality.

Gaussian Processes for Classification

- Objective: model posterior probabilities of the target variable for a new input
- Problem: we need to map values to interval $(0, 1)$
- Solution: use a Gaussian process together with a non-linear activation function (e.g., sigmoid)



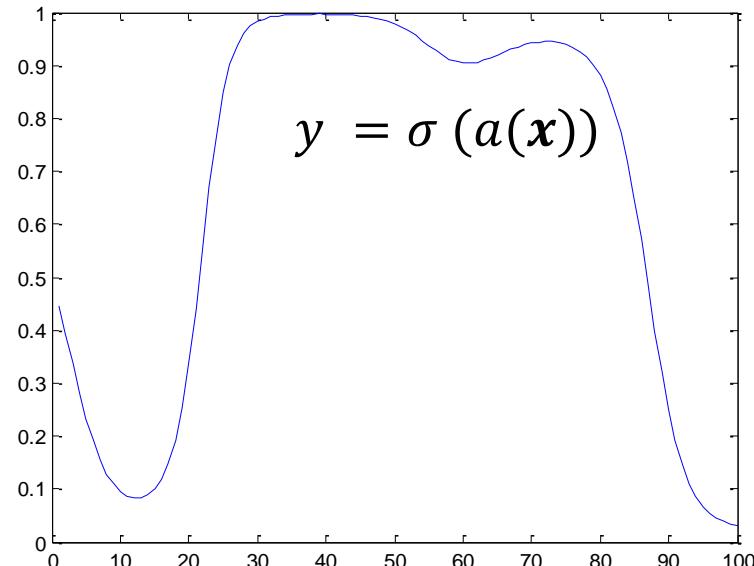
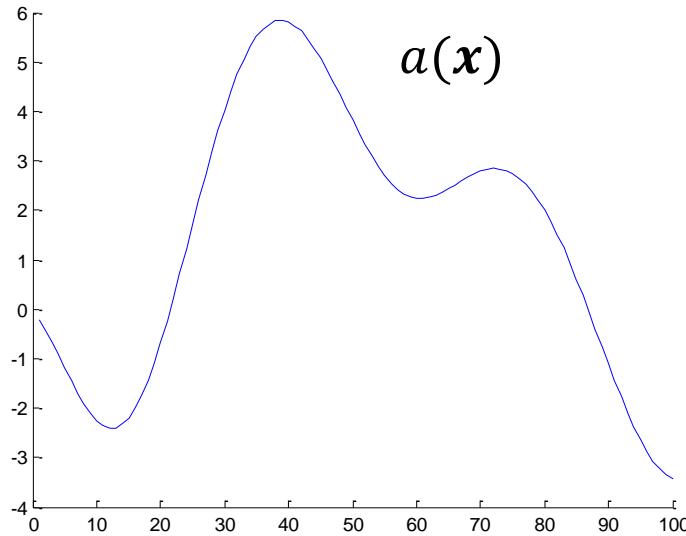
Sample from a Gaussian process over $a(x)$



MatLab Code

Gaussian Process for Classification

- Consider a two-class problem with target values $t \in \{0, 1\}$. Define a Gaussian process over a function $a(x)$ and transform a using the logistic sigmoid to $y = \sigma(a(x))$.
- This will give a non-Gaussian stochastic process over functions $y(x)$ where $y \in \{0, 1\}$.



[MatLab Code](#)

- The probability over the target t is a Bernoulli: $p(t | a) = \sigma(a)^t (1 - \sigma(a))^{1-t}$

Gaussian Process for Classification

- We denote the training set inputs by $\mathbf{x}_1, \dots, \mathbf{x}_N$ with corresponding observed target variables $\mathbf{t} = (t_1, \dots, t_N)^T$.
- We also consider a single test point \mathbf{x}_{N+1} with target value t_{N+1} .
- Our goal is to determine the predictive distribution $p(t_{N+1}|\mathbf{t})$, where we left the conditioning on the input variables implicit.
- Introduce a Gaussian process prior over the vector \mathbf{a}_{N+1} , which has components $a(\mathbf{x}_1), \dots, a(\mathbf{x}_{N+1})$.
- This defines a non-Gaussian process over \mathbf{t}_{N+1} , and by conditioning on \mathbf{t}_N we obtain the required predictive distribution $p(t_{N+1}|\mathbf{t})$.
- The Gaussian process prior for \mathbf{a}_{N+1} takes the form

$$p(\mathbf{a}_{N+1}) = \mathcal{N}(\mathbf{a}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$



Gaussian Process for Classification

- The covariance matrix does not include a noise term because we assume that all training data points are correctly labeled.
- However, for numerical reasons it is convenient to introduce a noise-like term governed by a parameter ν that ensures that the covariance matrix is positive definite. Thus the covariance matrix C_{N+1} has elements given by

$$C(x_n, x_m) = k(x_n, x_m) + \nu \delta_{nm}$$

- The kernel k is positive semidefinite and ν is fixed. $k(x, x')$ is governed by parameters θ that we will compute from the training data.
- For two class problems, it is sufficient to predict:

$$p(t_{N+1} = 1 | t_N)$$



Gaussian Process for Classification

- The predictive distribution is given as:

$$\begin{aligned} p(t_{N+1} = 1 | \mathbf{t}_N) &= \int p(t_{N+1} = 1, a_{N+1} | \mathbf{t}_N) da_{N+1} \\ &= \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} = \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \end{aligned}$$

- The integral is analytically intractable. Can be approximated with sampling (Neal)
- Integration can be analytical if $p(a_{N+1} | t_N)$ can be approximated with a Gaussian (Williams & Barber, Gibbs & MacKay, Gibbs)
- Approximations based on EP can also be used (Opper & Winther, Minka, Seeger)
 - Neal, R. M. (1997). [Monte Carlo implementation of Gaussian process models for Bayesian regression and classification](#). [Technical Report 9702](#), Department of Computer Statistics, University of Toronto.
 - Williams, C. K. I. and D. Barber (1998). [Bayesian classification with Gaussian processes](#). [IEEE Transactions on Pattern Analysis and Machine Intelligence](#) **20**, 1342–1351.
 - Gibbs, M. N. and D. J. C. MacKay (2000). [Variational Gaussian process classifiers](#). [IEEE Transactions on Neural Networks](#) **11**, 1458–1464.
 - Gibbs, M. N. (1997). [Bayesian Gaussian processes for regression and classification](#). Phd thesis, University of Cambridge.
 - Opper, M. and O. Winther (2000b). [Gaussian processes for classification](#). [Neural Computation](#) **12**(11), 2655–2684.
 - Minka, T. (2001b). [A family of approximate algorithms for Bayesian inference](#). Ph. D. thesis, MIT.
 - Seeger, M. (2003). [Bayesian Gaussian Process Models: PAC-Bayesian Generalization Error Bounds and Sparse Approximations](#). Ph. D. thesis, University of Edinburg.

Gaussian Process for Classification

- We need to compute $p(t_{N+1} = 1 | \mathbf{t}_N) = \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1}$
- One can compute the convolution of a Gaussian and a sigmoid function:

$$\int \sigma(a) \mathcal{N}(a|\mu, \sigma^2) da \simeq \sigma(\kappa(\sigma^2)\mu)$$

$$\kappa(\sigma^2) = \sqrt{1 + \frac{\pi\sigma^2}{8}}$$

We can thus try to approximate the posterior distribution $p(a_{N+1} | \mathbf{t}_N)$ as Gaussian

- How do we approximate the posterior?
 - Variational inference and make use of the local variational bound on the logistic sigmoid.
 - Expectation propagation
 - Laplace approximation



Laplace Approximation

- We can rewrite the posterior over a_{N+1} using Bayes' theorem:

$$\begin{aligned} p(a_{N+1} | \mathbf{t}_N) &= \int p(a_{N+1}, \mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1}, \mathbf{a}_N) p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1} / \mathbf{a}_N) p(\mathbf{a}_N) p(\mathbf{t}_N | \mathbf{a}_N) d\mathbf{a}_N \\ &= \int p(a_{N+1} / \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \end{aligned}$$

- We know how to compute mean and covariance for $p(\mathbf{a}_N | \mathbf{a}_N)$.

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}\left(\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}\right)$$

- It remains to find a Gaussian approximation only for $p(\mathbf{a}_N | \mathbf{t}_N)$
- This is done noting that $p(\mathbf{a}_N | \mathbf{t}_N) \propto p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N)$

Laplace Approximation

- We rewrote the posterior over a_{N+1} using Bayes' theorem:

$$p(a_{N+1} | \mathbf{t}_N) = \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N$$

where

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}\left(\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}\right)$$

- Also using $p(\mathbf{a}_N | \mathbf{t}_N) \propto p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N)$, $p(\mathbf{a}_N) = \mathcal{N}(0, \mathbf{C}_N)$ and

$$p(\mathbf{t}_N | \mathbf{a}_N) = \prod_{n=1}^N \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} = \prod_{n=1}^N e^{a_n t_n} \sigma(-a_n)$$

where in the last derivation we used: $\sigma(a_n) = \frac{1}{1 + e^{-a_n}}$ and

$$\sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} = \underbrace{(1 - \sigma(a_n))}_{\sigma(-a_n)} \left(\underbrace{\frac{\sigma(a_n)}{1 - \sigma(a_n)}}_{e^{a_n}} \right)^{t_n} = e^{a_n t_n} \sigma(-a_n)$$



Laplace Approximation

- The log of $p(\mathbf{a}_N / \mathbf{t}_N) \propto p(\mathbf{t}_N / \mathbf{a}_N)p(\mathbf{a}_N)$ which up to an additive constant is:

$$\begin{aligned} p(\mathbf{a}_N / \mathbf{t}_N) &\propto p(\mathbf{a}_N)p(\mathbf{t}_N / \mathbf{a}_N) \\ &\propto \mathcal{N}(0, \mathbf{C}_N) \prod_{n=1}^N e^{a_n t_n} \underbrace{\sigma(-a_n)}_{1/(1+e^{a_n})} \Rightarrow \end{aligned}$$

$$\begin{aligned} \Psi(\mathbf{a}_N) &= \ln p(\mathbf{a}_N / \mathbf{t}_N) = \ln p(\mathbf{t}_N / \mathbf{a}_N) + \ln p(\mathbf{a}_N) \\ &= -\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n}) \end{aligned}$$

- We will need to compute the 2nd derivative of this to come up with the Laplace approximation.

Laplace Approximation

$$\Psi(\mathbf{a}_N) = -\frac{1}{2}\mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n})$$

- The gradients are given as

$$\nabla \Psi(\mathbf{a}_N) = \mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N, \quad \boldsymbol{\sigma}_N = (\sigma(a_1), \dots, \sigma(a_N))^T$$

$$\nabla^2 \Psi(\mathbf{a}_N) = -\mathbf{W}_N - \mathbf{C}_N^{-1}, \quad \mathbf{W}_N = \text{diag}(\sigma(a_1)(1 - \sigma(a_1)), \dots, \sigma(a_N)(1 - \sigma(a_N)))$$

- We cannot find the mode by setting $\nabla \Psi(\mathbf{a}_N) = 0$ since $\boldsymbol{\sigma}_N$ depends nonlinearly on \mathbf{a}_N . For this we will use the IRLS algorithm.
- Note that $0 < \sigma(a_1)(1 - \sigma(a_1)) < 1/4 \Rightarrow \mathbf{W}_N$ is positive definite.
- The Hessian $A = -\nabla^2 \Psi(\mathbf{a}_N)$ is positive definite and thus the posterior $p(\mathbf{a}_N / \mathbf{t}_N)$ is log convex and has a single global mode.
- At the mode, the gradient vanishes giving: $\mathbf{a}_N^* = \mathbf{C}_N (\mathbf{t}_N - \boldsymbol{\sigma}_N)$

Laplace Approximation

- We can use the iterative reweighted least squares (IRLS) algorithm ($\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$) to update for \mathbf{a}_N (find the minimum of $-\ln p(\mathbf{a}_N / \mathbf{t}_N)$).
- The iterative update equation for \mathbf{a}_N is:

$$\mathbf{a}_N^{new} = \mathbf{C}_N (\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} (\mathbf{t}_n - \boldsymbol{\sigma}_n + \mathbf{W}_N \mathbf{a}_N)$$

Here we used

$$\begin{aligned}\mathbf{a}_N^{new} &= \mathbf{a}_N + (\mathbf{W}_N + \mathbf{C}_N^{-1})^{-1} (\mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N) = \mathbf{a}_N + \mathbf{C}_N (\mathbf{W}_N \mathbf{C}_N + \mathbf{I})^{-1} (\mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N) \\ &= \mathbf{C}_N (\mathbf{W}_N \mathbf{C}_N + \mathbf{I})^{-1} \{ (\mathbf{W}_N \mathbf{C}_N + \mathbf{I}) \mathbf{C}_N^{-1} \mathbf{a}_N + (\mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N) \} \\ &= \mathbf{C}_N (\mathbf{W}_N \mathbf{C}_N + \mathbf{I})^{-1} \{ \mathbf{W}_N \mathbf{a}_N + \mathbf{t}_N - \boldsymbol{\sigma}_N \}\end{aligned}$$

Laplace Approximation

- We can use the iterative reweighted least squares (IRLS) algorithm ($\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$) to update for \mathbf{a}_N .
- The iterative update equation for \mathbf{a}_N is:

$$\mathbf{a}_N^{new} = \mathbf{C}_N (\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} (\mathbf{t}_n - \boldsymbol{\sigma}_n + \mathbf{W}_N \mathbf{a}_N)$$

- The mode position \mathbf{a}_N^* and the Hessian matrix \mathbf{H} at this position $\mathbf{H} = -\nabla^2 \Psi(\mathbf{a}_N) = \mathbf{W}_N + \mathbf{C}_N^{-1}$ define the Gaussian approximation to $p(\mathbf{a}_N | \mathbf{t}_n)$

$$q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1})$$

In the Hessian note that \mathbf{W}_N is computed using \mathbf{a}_N^* .



Laplace Approximation

- Now we can go back to the formulas for linear Gaussian models and compute the integrals

$$p(a_{N+1} | \mathbf{t}_N) = \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N$$

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}\left(\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}\right)$$

$$p(\mathbf{a}_N | \mathbf{t}_N) \approx q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1}), \mathbf{H} = \mathbf{W}_N + \mathbf{C}_N^{-1}$$

from which we have:

$$p(a_{N+1} | \mathbf{t}_N) = \mathcal{N}\left(\mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N), c - \mathbf{k}^T (\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1} \mathbf{k}\right)$$

- Finally, we can also compute $p(t_{N+1} | \mathbf{t})$ (or the decision boundary $p(t_{N+1} | \mathbf{t}) = 0.5$ – in this case $\mu = \mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N) = 0$)

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1}$$

using $\int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da \approx \sigma(\kappa(\sigma^2)\mu)$, $\kappa(\sigma^2) = \sqrt{1 + \frac{\pi\sigma^2}{8}}$

Learning the Hyperparameters

- To determine the parameters θ of the covariance function, we can maximize the likelihood function:

$$p(\mathbf{t}_N | \theta) = \int p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N | \theta) d\mathbf{a}_N = \int p(\mathbf{t}_N, \mathbf{a}_N | \theta) d\mathbf{a}_N$$

- The integral is analytically intractable. The Laplace approximation can be applied again:
$$\ln p(\mathbf{t}_N, \mathbf{a}_N) = \Psi(\mathbf{a}_N) = -\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n})$$
$$\ln p(\mathbf{t}_N, \mathbf{a}_N) \approx \Psi(\mathbf{a}_N^*) - \frac{1}{2} (\mathbf{a}_N - \mathbf{a}_N^*)^T (\mathbf{W}_N + \mathbf{C}_N^{-1}) (\mathbf{a}_N - \mathbf{a}_N^*)$$
- Integrating in \mathbf{a}_N and using the normalization factor of a Gaussian gives:

$$\ln p(\mathbf{t}_N | \theta) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

where

$$\Psi(\mathbf{a}_N^*) = \ln p(\mathbf{a}_N^* | \theta) + \ln p(\mathbf{t}_N | \mathbf{a}_N^*)$$

- We now need an expression for the gradient of the log of $p(\mathbf{t}_N | \theta)$. Both \mathbf{a}_N^* and \mathbf{C}_N depend on θ

Learning the Hyperparameters

- To determine the parameters , maximize the likelihood:

$$\Psi(\mathbf{a}_N^*) = -\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N^* - \sum_{n=1}^N \ln(1 + e^{a_n^*})$$

$$\ln p(\mathbf{t}_N | \boldsymbol{\theta}) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

- Using $\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr}\left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x}\right)$ and $\frac{\partial \mathbf{A}^{-1}}{\partial x} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$, we obtain

the following terms of $\frac{\partial}{\partial \theta_j} \ln p(\mathbf{t}_N | \boldsymbol{\theta})$ with explicit dependence on $\boldsymbol{\theta}$:

$$\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left[\mathbf{C}_N^{-1} \left[\mathbf{I} - (\mathbf{W}_N + \mathbf{C}_N^{-1})^{-1} \mathbf{C}_N^{-1} \right] \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right] =$$

$$\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left[\mathbf{C}_N^{-1} \left[(\mathbf{C}_N \mathbf{W}_N + \mathbf{I})(\mathbf{C}_N \mathbf{W}_N + \mathbf{I})^{-1} - (\mathbf{C}_N \mathbf{W}_N + \mathbf{I})^{-1} \right] \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right]$$

$$\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left[(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{W}_N \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right]$$

Learning the Hyperparameters

- To determine the parameters, maximize the likelihood:

$$\Psi(\mathbf{a}_N^*) = -\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N^* - \sum_{n=1}^N \ln(1 + e^{a_n^*})$$

$$\ln p(\mathbf{t}_N | \boldsymbol{\theta}) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

- Using $\nabla \Psi(\mathbf{a}_N^*) = 0$, the term in $\frac{\partial}{\partial \theta_j} \ln p(\mathbf{t}_N | \boldsymbol{\theta})$ with explicit dependence on a_n^* is as follows:

$$-\frac{1}{2} \sum_{n=1}^N \frac{\partial \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}|}{\partial a_n^*} \frac{\partial a_n^*}{\partial \theta_j}$$

- Using $\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr}\left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x}\right)$ and $\mathbf{W}_N = \text{diag}(\sigma(a_i)(1 - \sigma(a_i)))$, it becomes $d\sigma / da_i = \sigma(a_i)(1 - \sigma(a_i))$

$$-\frac{1}{2} \sum_{n=1}^N \text{Tr} \left[(\mathbf{W}_N + \mathbf{C}_N^{-1})^{-1} \underbrace{\frac{\partial \mathbf{W}_N}{\partial \sigma(a_n)}}_{(1-2\sigma_n^*)} \underbrace{\frac{\partial \sigma(a_n)}{\partial a_n}}_{\sigma_n^*(1-\sigma_n^*)} \right] \frac{\partial a_n^*}{\partial \theta_j} = -\frac{1}{2} \sum_{n=1}^N \left[(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{C}_N \right]_{nn} \sigma_n^* (1 - \sigma_n^*) (1 - 2\sigma_n^*) \frac{\partial a_n^*}{\partial \theta_j}$$

Learning Hyperparameters

- We finally can compute $\partial \mathbf{a}_n^* / \partial \theta_j$ as follows : $\mathbf{a}_N^* = \mathbf{C}_N (\mathbf{t}_N - \boldsymbol{\sigma}_N) \Rightarrow$

$$\frac{\partial \mathbf{a}_N^*}{\partial \theta_j} = \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N) - \mathbf{C}_N \mathbf{W}_N \frac{\partial \mathbf{a}_N^*}{\partial \theta_j}$$

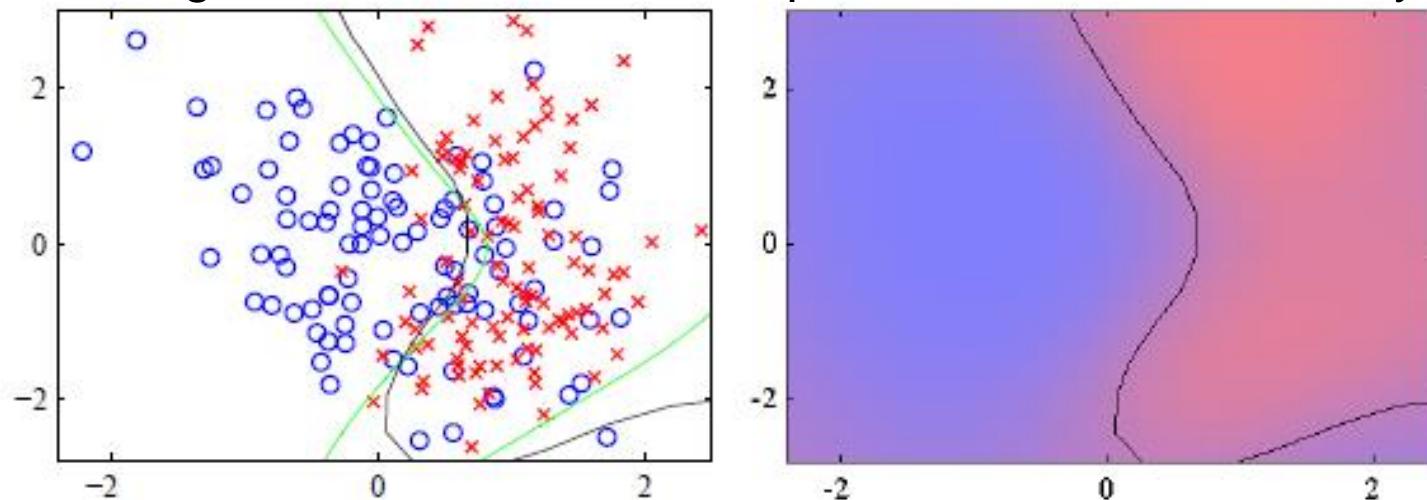
$\mathbf{W}_N = \text{diag}(\sigma(a_i)(1-\sigma(a_i)))$
 $d\sigma / da_i = \sigma(a_i)(1-\sigma(a_i))$

$$\frac{\partial \mathbf{a}_N^*}{\partial \theta_j} = (\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N)$$

- We now have all the terms in the gradient $\frac{\partial}{\partial \theta_j} \ln p(\mathbf{t}_N | \theta)$ and can employ gradient optimization techniques.

Binary Classifier Based on Gaussian Process

- We now have a binary classifier based on a Gaussian process. On left the data together with the optimal decision boundary from the true distribution in green, and the decision boundary from the Gaussian process classifier in black.
- On right the predicted posterior probability for the blue and red classes together with the Gaussian process decision boundary.



Requires installing [NetLab](#)
and [setting the appropriate path](#)

[MatLab Code](#)

- Williams, C. K. I. and D. Barber (1998). [Bayesian classification with Gaussian processes](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence* **20**, 1342–1351.



Connections to Neural Networks

□ Neural Networks

- The range of representable functions is governed by the number M of hidden units
- Within the maximum likelihood framework, they overfit as M comes close to the number of training samples

□ Bayesian Neural Networks

- The prior over w in conjunction with the network function $f(x, w)$ produces a prior distribution over functions from $y(x)$
- The distribution of functions will tend to a GP as $M \rightarrow \infty$.
 - In this limit the output variables of the neural network become independent.
 - However, the property that the outputs share hidden units and thus all of them affecting the values of the weights is lost in this limit.

▪ Neal, R. M. (1996). *Bayesian Learning for Neural Networks*. Springer. Lecture Notes in Statistics 118.



Connections to Neural Networks

- There are explicit forms for the covariance in the case of probit and Gaussian hidden unit activation functions.
 - These kernel functions $k(x, x')$ are nonstationary, i.e. not a function of $x - x'$.
 - This is a consequence of the Gaussian weight prior being centered on zero which breaks translation invariance in weight space.
- By working directly with the covariance function we have implicitly marginalized over the distribution of weights.
- If the weight prior is governed by hyperparameters, then their values will determine the length scales of the distribution over functions.

▪ [Williams, C. K.](#) I. (1998). [Computation with infinite neural networks](#). [Neural Computation 10\(5\), 1203–1216](#).

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)

